



# Agentic Runtime Reconfiguration of Architectural Patterns in Federated Learning<sup>☆</sup>

Ivan Compagnucci<sup>a</sup>, Qinghua Lu<sup>b</sup>, Catia Trubiani<sup>a</sup>

<sup>a</sup> Gran Sasso Science Institute, Via Michele Iacobucci 2, 67100, L'Aquila, Italy

<sup>b</sup> Data61, CSIRO, Level 5/13 Garden St, Eveleigh, Australia

## ARTICLE INFO

### Keywords:

Architectural Patterns  
Federated Learning  
AI Agents  
Large Language Models

## ABSTRACT

Federated Learning (FL) allows multiple clients to collaboratively train a global model without sharing their local data, thereby preserving privacy. However, architecting FL systems involves complex architectural design choices that depend on both intrinsic system characteristics (e.g., the number of available clients) and evolving runtime conditions, such as network instability or new data arriving between training rounds. While architectural patterns can address these concerns, existing approaches typically keep them static throughout the entire FL execution. This static activation poses a significant limitation, as a pattern that provides benefits in early training rounds may later introduce performance drawbacks as operational conditions evolve. For instance, a pattern intended to improve model accuracy might eventually introduce more time and computational overhead if it is not deactivated when no longer needed. This paper introduces an agentic approach that treats architectural patterns as runtime-reconfigurable entities. We extend the standard FL workflow with a Multiple AI-Agent layer that leverages Large Language Models (LLMs) to dynamically (de)activate patterns after each training round. These decisions are driven by real-time performance metrics and evolving system characteristics. To coordinate the agents, we implement and compare three distinct coordination strategies: Voting-based, Role-based, and Debate-based, evaluating them against three baselines. Experiments are conducted by emulating clients with heterogeneous characteristics and varying system operational conditions to test the effectiveness of our approach. Results examine different datasets and show that, by adaptively (de)selecting architectural patterns across FL rounds, agentic approaches improve learning accuracy and can reduce execution time, at the cost of a moderate overhead.

## 1. Introduction

Federated Learning (FL) has emerged as a novel paradigm in distributed machine learning, enabling the orchestration of global model training without sharing raw data (Hongyi Zhang et al., 2020; McMahan et al., 2017; Kairouz et al., 2021). By decentralizing the learning process, FL allows multiple clients to collaboratively train a model on their local data, addressing critical privacy concerns. This way, each client shares only the trained model weights with a centralized server, which then aggregates them into a global model using information fusion algorithms (e.g. FedAvg McMahan et al., 2017).

In the software architecture community, designing FL systems is recognized as a primary challenge, with *performance optimization* being a central concern (Kairouz et al., 2021; Hongyi Zhang et al., 2020). This is coherent with studies on distributed computing environments, which show that dynamic workloads and resource heterogeneity can affect execution efficiency (Menon et al., 2024; Verbraeken et al., 2021). While

significant research in FL has focused on refining training algorithms and global model aggregation techniques, the role of the underlying architecture has received less attention. A seminal contribution in this context is the FLRA reference architecture (Lo et al., 2021), which promotes the use of architectural patterns to encapsulate best practices for recurring problems. Building upon this reference architecture, the same authors further identified a catalog of 15 architectural patterns specifically tailored for FL systems (Lo et al., 2022). Concrete examples of such patterns include: (i) the *Client Selector*, which restricts training participation to clients meeting specific criteria; (ii) the *Heterogeneous Data Handler*, which employs data augmentation to mitigate training data heterogeneity; and (iii) the *Message Compressor*, which compresses the size of client-server messages to minimize communication latency.

Traditionally, architectural patterns in FL are implemented statically at design time, meaning that the list of active patterns remains

<sup>☆</sup> This article is part of a Special issue entitled: 'AI for SA' published in The Journal of Systems & Software.

\* Corresponding author.

E-mail addresses: [ivan.compagnucci@gssi.it](mailto:ivan.compagnucci@gssi.it) (I. Compagnucci), [qinghua.lu@data61.csiro.au](mailto:qinghua.lu@data61.csiro.au) (Q. Lu), [catia.trubiani@gssi.it](mailto:catia.trubiani@gssi.it) (C. Trubiani).

fixed throughout the entire training process (Chen Zhang et al., 2021; Lo et al., 2022). Our previous works (Compagnucci et al., 2026b, 2025) show that these patterns exhibit both benefits and drawbacks depending on system settings (e.g., the number of clients and the efficiency of the communication network). However, since these settings can change at runtime (e.g., the communication network may become unstable), a static architectural pattern setup fails to maximize the effectiveness of the FL system. In fact, a pattern that is beneficial in the initial rounds may introduce drawbacks later. To provide a practical example, consider a *Client Selector* pattern whose selection criterion is based on computational power. In the early rounds, including all clients may be ideal to maximize data coverage and model generalization. If network conditions deteriorate in later rounds, dynamically activating this pattern to exclude low-power clients can prevent bottlenecks and significantly reduce the total time required to complete the FL execution.

Based on these premises, we envision architectural patterns as *runtime-reconfigurable entities* that can be adaptively (de)activated at each round of the FL execution. This implies monitoring the actual system's configuration and deciding which patterns contribute to enact the most effective design alternative at runtime. However, managing such complex decisions (i.e., which pattern (de)activate and why) poses significant challenges. The decision-making process must account for unpredictable factors, such as fluctuating network latency and varying data availability, which can change throughout the training process. Therefore, this requires a continuous, round-by-round evaluation that correlates the actual system settings with performance metrics gathered from previous rounds. This need aligns with a growing trend in the software architecture community, where Large Language Models (LLMs) are increasingly used to support complex architectural decisions (Soliman and Keim, 2025; Li et al., 2024; Pace et al., 2024, 2025; Dhar et al., 2024a,b; Compagnucci and Trubiani, 2025). Beyond decision support, LLM-based agents can be integrated into software systems to monitor runtime behavior, learn from historical outcomes, and suggest real-time interventions. This integration is particularly valuable in the FL context, enabling systems to autonomously self-optimize and adapt to evolving operational contexts (Bass et al., 2025; Vaidhyathan and Muccini, 2025).

This work presents an agentic approach for the dynamic runtime management of architectural patterns in FL systems. To achieve this, we first instantiate the FLRA reference architecture by integrating three specific patterns from Lo et al. (2022), namely Client Selector, Heterogeneous Data Handler, and Message Compressor. We then extend this architecture by introducing a Multiple AI-Agent layer. At the conclusion of each FL round, this layer evaluates the current system settings and performance metrics collected so far to decide the configuration of active patterns expected to be most effective. To govern this decision-making process, we implement and compare three coordination strategies proposed in the AgentOps catalog from Liu et al. (2025): (i) Voting-based, where agents submit preferences and a coordinator takes the final decision; (ii) Role-based, where agents are assigned specific roles, such as specialist for a specific pattern, to reach a conclusion; and (iii) Debate-based, where agents provide arguments to motivate their opinions and interact to converge toward a consensus. The effectiveness of our approach is validated through an extensive experimental campaign designed to grasp the inherent complexity and heterogeneity of FL systems. Specifically, we conduct experiments in an FL environment characterized by high system and data heterogeneity, where clients exhibit diverse computational capacities, skewed data distributions, and different data-inflow regimes, all under stochastic network instability. Results show that architectural decisions, i.e., the choice of whether to activate or deactivate a pattern, can be transformed into effective runtime control mechanisms. By introducing an agentic approach to coordinate the activation of architectural patterns, the system effectively turns static design choices

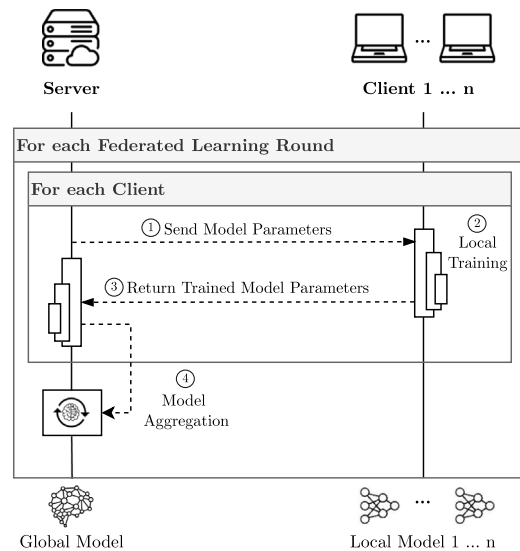


Fig. 1. Federated learning workflow.

into active entities, enabling the architecture to autonomously adapt and optimize its efficiency in real-time.

In summary, the main contributions of this paper are: (i) we extend the Federated Learning Reference Architecture (FLRA) by introducing a Multiple AI-Agents layer that treats architectural patterns as *runtime configurable entities*. The system dynamically (de)activates patterns by reasoning over evolving operational conditions and real-time performance metrics; (ii) we implement three distinct agentic coordination strategies to automate the selection of three actionable architectural patterns; (iii) we conduct an extensive experimental campaign that provides quantitative evidence of the Multiple AI-Agents layer effectiveness in representative FL scenarios. To ensure reproducibility, all artifacts are publicly available (Compagnucci et al., 2026a).

## 2. Preliminaries

### 2.1. Federated learning in a nutshell

In traditional Machine Learning (ML), training data is typically centralized on cloud servers to enable the development of a single predictive model (Amershi et al., 2019). However, rising privacy concerns and regulatory constraints have reduced the willingness to share sensitive information (Kairouz et al., 2021; Sánchez et al., 2024), resulting in data silos where data are isolated for protection (Yang et al., 2019). Federated Learning (FL) offers a privacy-preserving alternative by enabling collaborative model training across distributed data sources without transferring raw data (Kairouz et al., 2021). Due to these advantages, FL has experienced rapid adoption in recent years, with successful applications across various domains, including healthcare, finance, and industry (Hongyi Zhang et al., 2020).

The main steps of the FL workflow are illustrated in Fig. 1. The workflow starts with a central server broadcasting the initial global model parameters (e.g., model weights) to all participating clients ①. Each client then trains the model locally using its private data ② and returns only the updated model weights to the server ③. The server aggregates these updates to obtain a new global model ④, which is subsequently redistributed to the clients for the next training round. This cycle constitutes a single FL round and is repeated for a fixed number of rounds or until a predefined stopping criterion is met (e.g., convergence of the global model, stabilization of validation performance).

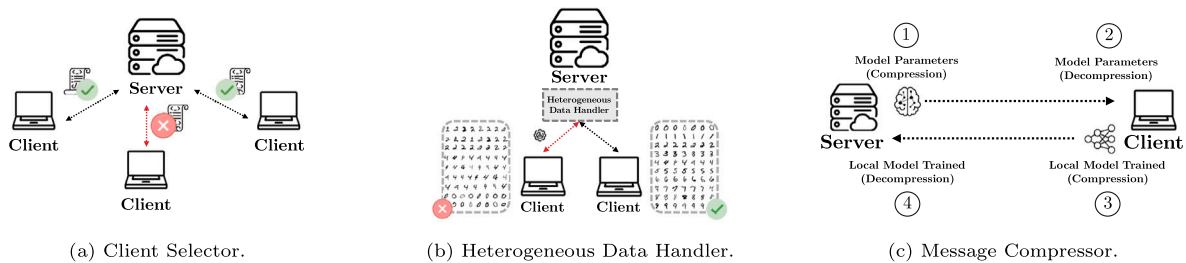


Fig. 2. Architectural patterns candidates for selection.

## 2.2. Architectural patterns in federated learning systems

Architectural patterns offer reusable solutions for common design challenges in complex systems (Richards, 2015). Lo et al. (2022) propose a set of patterns tailored to FL, addressing macro areas such as client management, model management, and model training, covering these system dimensions with architectural solutions. In our study, we consider different dimensions, focusing on clients, data, and communication, which are targeted by the following patterns: the Client Selector, the Heterogeneous Data Handler, and the Message Compressor (Lo et al., 2022), as shown in Fig. 2. It is worth recalling that patterns are evaluated (and possibly selected) at each round of FL execution, thereby accounting for dynamic changes during operation.

### 2.2.1. Client selector

The Client Selector is depicted in Fig. 2(a). It determines the subset of clients that will participate in the learning process according to predefined criteria (Lo et al., 2022). Selection criteria can be: (i) *data-centric* (e.g., dataset size, heterogeneity, quality), (ii) *resource-centric* (e.g. computation and network capacity), and (iii) *performance-centric* (i.e., recent contribution to the global model) (Briggs et al., 2020; Lo et al., 2022). Our current implementation includes the *resource-centric* selection criterion: clients are excluded if the number of physical CPU cores per client is below a predefined threshold.

**Performance Implications.** Client Selector can reduce both training and total round time, by excluding clients with low computational power, thereby mitigating bottlenecks during global model aggregation (Lo et al., 2022; Compagnucci et al., 2026b; Vu et al., 2021; Compagnucci et al., 2025). On the other hand, it may slow model convergence or degrade model accuracy, because excluding clients from training prevents the model from incorporating weight updates computed from their local data.

### 2.2.2. Heterogeneous data handler

Fig. 2(b) depicts the goal of the Heterogeneous Data Handler, which is to mitigate issues arising from non-IID client data (e.g., local datasets with different label proportions or feature distributions) that may negatively affect the global model’s accuracy. In our implementation, the pattern applies *data augmentation* to generate synthetic samples and increase the diversity and size of the local dataset using Generative Adversarial Networks (GANs) (Goodfellow et al., 2014; Yu et al., 2017; Zhang et al., 2017). When active, a generator produces synthetic samples while a discriminator learns to distinguish them from real ones; as training progresses, the generator improves and is used to create class-conditional samples that populate underrepresented classes. Non-IID conditions are detected through the Jensen–Shannon Divergence (JSD) metric (Hu et al., 2025) computed locally on each client and shared with the server as a single scalar  $JSD_i^r = JSD(D_i^r) \in [0, 1]$  for round  $r$ , where  $D_i^r$  is the client dataset and  $JSD_i^r = 0$  denotes a perfectly balanced distribution. This supports server-side decision making without exposing raw data or dataset statistics (Hu et al., 2025).

**Performance Implications.** This pattern improves model accuracy by mitigating the effects of non-IID data through augmentation. However,

it introduces additional overhead due to the extra steps required by the GAN augmentation mechanism. In our setting, this overhead is included in the training time because, when enabled, GAN-based augmentation runs as an additional sub-phase before the client’s local training (Lo et al., 2022; Compagnucci et al., 2026b).

### 2.2.3. Message compressor

The Message Compressor aims to reduce the communication time overhead by compressing the model parameters exchanged between the server and clients. Fig. 2(c) illustrates the four phases (Lo et al., 2022): ① the server compresses the model weights and sends them to clients; ② clients decompress the received parameters for local training; ③ after training, clients compress their updated parameters and return them to the server; ④ the server decompresses incoming updates for aggregation. Weight compression can reduce communication overhead when the client–server payload exchanged is large, e.g., models with a large number of parameters (Lo et al., 2022; Hongyi Zhang et al., 2020). Compression yields benefits only if the data exchanged is sufficiently large; otherwise, the compression and decompression overhead outweigh the benefits (Compagnucci et al., 2025).

**Performance Implications.** By reducing the volume of data transmitted, this pattern can decrease Client–Server communication time and save bandwidth, which is particularly beneficial for bandwidth-constrained clients (Lo et al., 2022). However, the extra computation required for compression and decompression may outweigh the communication-time savings, especially when the exchanged payload is small (Hongyi Zhang et al., 2020; Lo et al., 2022; Compagnucci et al., 2025).

## 3. Related work

In this section we review three complementary lines of work linked to our contribution: (i) the use of LLMs to support software-architecture tasks and architectural decision making, (ii) architectural solutions for FL, including reference architectures and catalogs of FL architectural patterns, and (iii) dynamic and self-adaptive FL approaches where the *FL system* adapts its architecture or behavior on the fly. Building on these three macro areas, we then position our contribution with respect to prior work by clarifying the novelty of our approach.

### 3.1. LLMs for software architecture design

Software engineers increasingly adopt LLMs for a variety of development activities, and software architects have begun to investigate their suitability for supporting architectural design decisions. LLMs are increasingly used to support architectural design decisions (Pace et al., 2024; Soliman and Keim, 2025; Li et al., 2024; Arun et al., 2025). Soliman and Keim (2025) evaluate seven LLMs on questions about the architectural knowledge of existing systems, showing that LLMs encode useful architectural information but suffer from precision and reliability issues. Dhar et al. (2024a) conduct an empirical study on LLMs generating architectural design decisions as architecture decision records. Preliminary analyses suggest that such LLM-generated decisions can support software architects, although further research is needed. Pace

et al. (2024) introduce an LLM-based assistant that guides novice architects through interactive decision-making workflows that combine architectural knowledge and quality attributes. Dhar et al. (2024b) study whether LLMs can generate Architecture Decision Records and later propose DRAFT, a retrieval-augmented and few-shot fine-tuned approach for more effective and efficient ADR drafting. Arun et al. (2025) investigate the effectiveness of LLMs for generating architectural components by comparing LLM-produced code with human-developed implementations.

### 3.2. Architectural solutions in federated learning

Prior work has proposed several architectural solutions for FL systems. FLRA (Lo et al., 2021) defines a pattern-oriented reference architecture spanning the main FL phases, from job creation to monitoring, and provides the foundation for our Multiple AI-Agents layer. Baresi et al. (2025) present a requirement-driven reference architecture that integrates node selection and configuration strategies and empirically compares alternative configurations under varying conditions. Hongyi Zhang et al. (2020) analyze four FL system architectures (centralized, hierarchical, regional, decentralized) and compare them in terms of communication overhead, convergence speed, and scalability. Lo et al. (2022) complement these contributions with a catalog of FL architectural patterns (e.g., the Client Selector adopted in this work), each providing actionable design strategies.

Many research works benchmark the performance of FL systems. FedScale (Lai et al., 2022) offers realistic datasets, a scalable runtime, and high-level APIs to support and benchmark FL experiments. Li et al. (2020) survey FL optimization issues, highlighting the impact of device constraints, non-IID data, and privacy. Client selection strategies have been systematically analyzed (Fu et al., 2023; Mayhoub and M. Shami, 2024), addressing data heterogeneity, hardware limitations, and scheduling. Compagnucci et al. (2026b) empirically evaluate four FL architectural patterns and their combinations from Lo et al. (2022), quantifying the trade-offs between computational overhead and resulting performance gains.

### 3.3. Architectural adaptation in federated learning

Dynamic FL focuses on developing methods or strategies to improve the efficiency of FL executions at each round (Rizk et al., 2020). However, these methods are typically driven by runtime observations (e.g., accuracy trends, resource availability, and data heterogeneity), while keeping the underlying system architecture fixed (Rizk et al., 2020). As a result, architectural variability is not treated as a first-class adaptation mechanism, and the architecture itself is not leveraged as a runtime optimization knob. Baresi et al. (2021) frame FL as a self-adaptation problem: given a target accuracy for the next round, the system estimates the number of local epochs based on clients' resource constraints and the accuracy observed in the previous two rounds. Wang et al. (2019) introduce a control scheme that, under a fixed resource budget, tunes the balance between local computation and global aggregation across rounds to reduce training loss. Li et al. (2021) propose a server-side approach that exploits information from past executions to predict, for each client, the workload it can realistically sustain in subsequent rounds. Jie Zhang et al. (2021) mitigate the degradation caused by non-IID data by adjusting each client's batch size at every round, tailoring local training to the observed heterogeneity. Ilhan et al. (2023) address model downsizing by adding early-exit classifiers, thereby improving training cost-effectiveness for clients with limited resources. Singh and Adhikari (2025) address the problem of data heterogeneity via an FL strategy that applies adaptive masking to unlabeled data, thereby reducing prediction entropy and increasing confidence. These contributions confirm the usefulness of self-adaptation in FL. However, they mostly target training-level levers (e.g., the number of local epochs, batch size, aggregation frequency) or

the learning model itself (e.g., architectural variants or hyperparameter settings), rather than the FL solution's system architecture. In contrast, our work applies self-adaptation to architectural decisions by selecting and dynamically toggling FL architectural patterns, which provides a distinct and complementary optimization dimension for FL systems.

### 3.4. Research gap and positioning

Existing research on FL architectures provides reference models and catalogs of architectural patterns, but these patterns are typically treated as static design-time decisions (Lo et al., 2022, 2021). Once selected, they remain statically configured throughout the training process, even though system conditions (e.g., client heterogeneity, network instability, and evolving data availability) may change across rounds. At the same time, dynamic FL approaches focus primarily on adapting training-level parameters (e.g., local epochs, aggregation frequency, batch size), while leaving the architectural structure unchanged. As a result, architectural variability has not been systematically exploited as a runtime optimization mechanism.

Our work addresses this gap by elevating architectural pattern activation to a first-class runtime architectural decision. We leverage LLM-based agents not merely as advisory tools, but as decision-making components that reason over system settings and historical evaluation metrics to determine, at each round, the most suitable architectural setup. By dynamically reconfiguring the system architecture itself, rather than only tuning training parameters, we introduce a complementary adaptation dimension for FL systems. This positions our contribution at the intersection of FL architectural design and agent-based runtime adaptation, enabling architectural decisions to be treated as systematic and context-aware optimization entities.

## 4. Our approach

In this section, we present the rationale for our work and describe the architectural design of our approach, which instantiates FLRA (Lo et al., 2021) and augments it with a Multi AI-Agents layer. We then present three coordination strategies implemented to manage agent decisions, explaining how this layer integrates into the standard FL workflow.

### 4.1. Motivation

Let us consider the goal of pursuing the *performance optimization* in FL systems. A traditional static configuration of architectural patterns acts like a fixed checklist, applying the same setup each FL round, regardless of context. In contrast, an AI-based agent may act like an *intelligent system architect*, analyzing past and actual performance metrics and dynamically selecting architectural patterns for the next FL round. The operational logic of this agentic layer can be grounded in the MAPE-K (Monitor-Analyze-Plan-Execute over Knowledge) reference model for self-adaptive systems (Kephart and Chess, 2003; Brun et al., 2009). From this perspective, the agent realizes a structured feedback loop that *monitors* the FL environment, *analyzes* performance and system metrics, *plans* the possible architectural alternatives, and *executes* the selection of the most appropriate pattern(s) for the subsequent round. By placing AI agents at the core of this loop, the framework leverages their reasoning capabilities as the *Knowledge* component, thus the system can autonomously manage complex trade-offs and adapt to changing operational conditions, thereby continuously monitoring the FL system.

Consider the use case reported by Dayan et al. (2021), where a group of researchers deployed a FL system for the COVID-19 emergency, where 20 hospitals collaboratively trained a classification model without moving patient data, using only the first emergency-department measurements available at presentation (i.e., vital signs, laboratory values, and the initial chest X-ray) to predict patients'

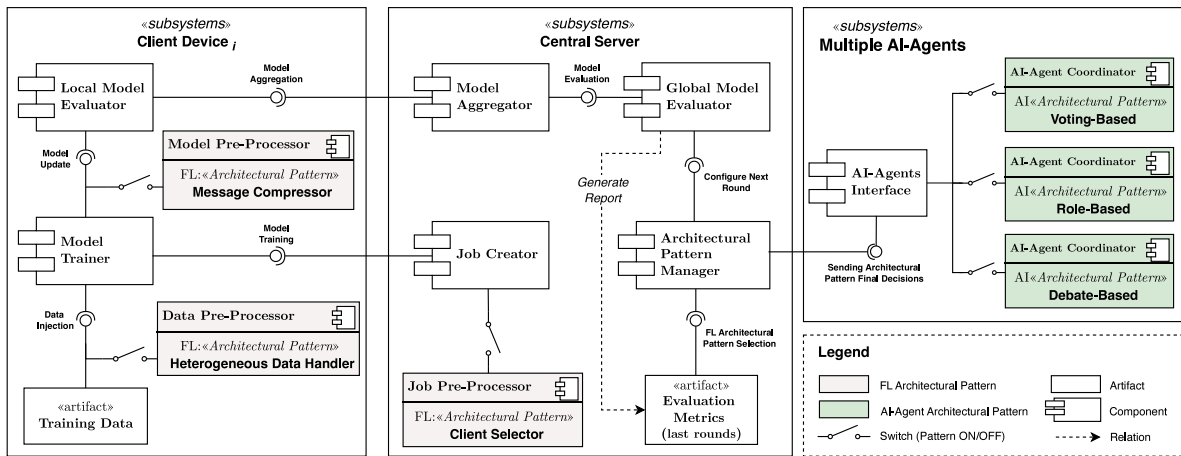


Fig. 3. Component diagram representing the FLRA extended with the multiple AI-agents layer.

oxygen needs at 24 and 72 h. In such settings, early rounds may benefit from selecting a larger set of hospitals (*Client Selector*) to maximize data coverage and train the global model with a larger dataset. As training progresses, however, including every hospital in every round can make the overall process unnecessarily slow: clients with limited computational resources or poor connectivity may consistently take longer to complete local training, becoming stragglers that delay aggregation and turn each round into a bottleneck. In later rounds, the client selector can therefore prioritize clients with higher computational capacity, reducing round duration and accelerating convergence. Likewise, when the local datasets follow different data distributions across hospitals, enabling a *Heterogeneous Data Handler* can mitigate non-IID effects (e.g., skewed case severity across hospitals) and improve generalization; and when network latency becomes a bottleneck, activating a *Message Compressor* can reduce communication time by exchanging compressed updates. The key advantage is that these patterns can be (de)activated on the fly, making the FL process dynamic and context-aware, thereby improving overall system performance.

#### 4.2. Extending the FLRA with a multiple AI-agents layer

Our approach extends the Federated Learning Reference Architecture (FLRA) (Lo et al., 2021) by introducing a Multiple AI-Agents layer. The rationale for choosing this reference architecture is that, to the best of our knowledge, it is the only one that enables modular integration of architectural patterns while preserving the standard FL workflow (see Fig. 1). Moreover, FLRA is grounded in a systematic literature review and qualitatively validated using evidence from both the literature and industrial implementations (Lo et al., 2021). FLRA's explicit modularization of architectural pattern components enables runtime reconfiguration, allowing us to (de)activate and combine patterns on the fly without affecting the rest of the system. This choice also supports our optimization goals, as the architectural patterns implemented in FLRA affect system architecture (e.g., client participation, local training configuration, and communication policies), which, in turn, influence system evaluation metrics such as accuracy, training time, and communication time.

Fig. 3 depicts a component diagram representing the FLRA architecture extended with the Multiple AI-Agents layer. It consists of three components: the *Central Server*, the *i*-th *Client Device* (i.e., multiple instances may be present; only one is illustrated for simplicity), and the *Multiple AI-Agents* layer. Note that FLRA includes the FL architectural patterns analyzed in this study (i.e., *Client Selector*, *Message Compressor*, and *Heterogeneous Data Handler*), which represent design strategies of the FL system itself. The *Multiple AI-Agents* layer includes three coordination strategies (i.e., *Voting-based*, *Role-based*,

and *Debate-based*) to manage and coordinate the AI agents' decision-making processes. Each architectural pattern and coordination strategy is modeled as a configurable component that can be activated (ON) or deactivated (OFF). Importantly, while all patterns can be switched OFF, at least one coordination strategy must be enabled to ensure agentic decision-making.

*Central server.* It initializes the ML task, coordinates federated training, collects client updates, aggregates the global model, and computes core metrics. The server includes the *ARCHITECTURAL PATTERN MANAGER* to determine which patterns to enable, and applies this decision by toggling the pattern switches, thereby configuring the patterns enabled for the next FL round. Specifically, the component manager receives a binary array from the agentic layer, where each value indicates whether a target pattern is active.

*Client device<sub>i</sub>.* Each client prepares its local data and model, trains and evaluates locally, and sends its update to the server. At the beginning of every FL round, the *ARCHITECTURAL PATTERN MANAGER*, guided by the *Multiple AI-Agents* layer, broadcasts a binary array to each client specifying which patterns to activate or not. Patterns apply to three stages of the workflow: (i) selecting which clients join the round (*Client Selector*); (ii) augmenting client training data (*Heterogeneous Data Handler*); and (iii) compressing client-server messages (*Message Compressor*).

*Multiple AI-agents.* This layer consists of a set of AI agents that leverage LLMs to analyze the current FL system configuration and the evaluation metrics collected so far, identifying the pattern combination they deem optimal for the next round. Note that both architecture pattern activation and deactivation do not require additional client-side infrastructure; clients simply execute the corresponding modules as part of the standard FL workflow. This avoids compatibility issues across heterogeneous clients and prevents additional runtime overhead beyond the execution of the selected architectural pattern modules.

The agentic recommendation process employs three coordination strategies drawn from the AgentOps catalog (Liu et al., 2025), specifically: *Voting-based*, *Role-based*, and *Debate-based Cooperation*. Additional details on these coordination strategies and their implementation are provided hereafter.

#### 4.3. Multiple AI-agents coordination strategies

##### 4.3.1. Voting-based cooperation

This strategy introduces a decision mechanism in which multiple AI agents independently form opinions and then cast votes to a centralized coordinator agent, which then elaborates on and delivers a verdict (Liu et al., 2025). Typical voting rules include (i) *democratic majority* or (ii) *weighted voting*, where weights may reflect agent roles or reliability.

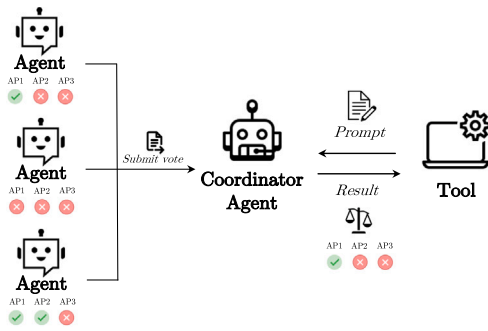


Fig. 4. Voting-based coordination strategy.

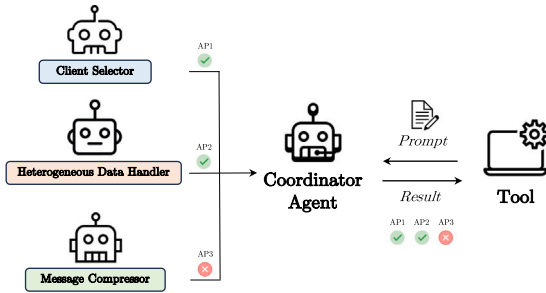


Fig. 5. Role-based coordination strategy.

As depicted in Fig. 4, a central coordinator agent receives a prompt incorporating a task from a tool. It then broadcasts the request to a set of agents, who reason locally and return their votes. The coordinator evaluates the votes and, according to the chosen voting rule, sends the final decision, i.e., the set of patterns to be activated in the next FL round. This strategy preserves fairness and accountability by logging each agent’s vote, allowing for the aggregation of collective intelligence to enhance decision quality.

*Our Implementation.* We implement the Voting-based coordination strategy by deploying multiple AI agents alongside a coordinator agent. At the end of each FL round, each agent, together with a textual explanation, can express a binary preference (ON or OFF) for each pattern, which will then be compared with the preferences expressed by other agents on the same pattern. The coordinator collects these choices and decides whether to set the pattern to ON or OFF based on the democratic majority, i.e., the pattern endorsed by receiving >50% of votes from agents is selected.

#### 4.3.2. Role-based cooperation

This strategy assigns explicit roles to the agents, each with specific responsibilities (Liu et al., 2025). Typical roles include (i) *coordinator* to orchestrate other agents; (ii) a *specialist*, which performs predefined and tailored tasks; and (iii) *evaluator/auditor* to assess outcomes and compliance. As depicted in Fig. 5, a central coordinator agent receives a task, applies the role policy, and dispatches tasks to agents according to their roles; agents act within their corresponding role boundaries and return outputs accordingly. The coordinator then merges and organizes the results. This strategy improves accountability and safety by enforcing least privilege, separating concerns, and making each action attributable to an agent, which simplifies auditing (Liu et al., 2025).

*Our Implementation.* We implement the Role-based coordination strategy by defining two main agent roles: the *Coordinator* and the *Pattern Architect*. Each Pattern Architect is “specialized” in a specific architectural pattern under analysis (i.e., Client Selector, Heterogeneous Data Handler, and Message Compressor) and implemented as an AI agent. The task of each Pattern Architect is to recommend whether the corresponding pattern should be selected in the next round, and

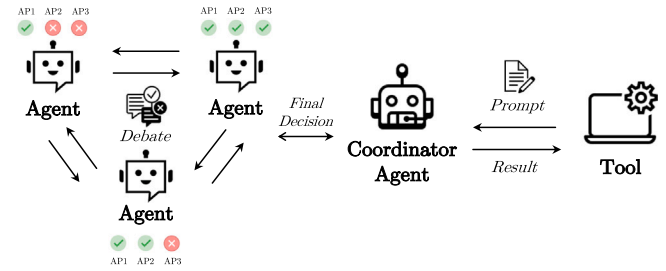


Fig. 6. Debate-based coordination strategy.

justify the decision with a rationale. The Coordinator then aggregates all recommendations and applies the resulting configuration, that is, the combination of ON/OFF pattern states for the next FL round.

#### 4.3.3. Debate-based cooperation

The Debate-based cooperation strategy lets agents present proposals and challenge each other’s reasons under a coordinator agent (Liu et al., 2025). It is useful when trade-offs are complex and evidence is uncertain, as arguments must be explicit, comparable, and auditable. Typical rules include time-limited turns, a single critique round, and scoring mechanisms to select the winning proposal. As depicted in Fig. 6, the coordinator proposes a task along with its corresponding context. Based on this information, agents submit a proposal accompanied by a rationale. Subsequently, agents criticize and may revise their proposals. After this revision, the coordinator applies the scoring rule to select the best-supported option and outputs the decision. This strategy enhances transparency and decision quality by requiring agents to provide reasons tied to evidence, thereby reducing single-heuristic bias and making the outcome easier to audit (Liu et al., 2025).

*Our Implementation.* The Debate-based coordination strategy is implemented by deploying multiple AI agents and a coordinator agent. At the end of each FL round, the coordinator triggers a debate among the AI agents. Each agent first drafts its ON/OFF proposal for all patterns with a brief rationale. Then, agents exchange their positions and iteratively revise both proposal and critique in reply to other agents’ proposals. A lightweight shared debate memory records arguments and refinements. The debate stops when agents converge on the same proposal or when a maximum of  $n$  turns is reached. As suggested in Liu et al. (2025), in our experimentation, we set a maximum number of debate turns (i.e.,  $n = 5$ ) to avoid the overhead introduced due to the high reasoning time. If no consensus emerges, the current configuration is carried over to the next round, with convergence defined as stable proposals or a strict majority agreeing on each pattern.

#### 4.4. Integrating the multiple AI-agents layer into the federated learning workflow

Fig. 7 depicts the extended FL workflow, introducing the “Architectural Pattern Selection” phase ⑤. After the “Model Aggregation” step ④, the server maintains a **system configuration file** reflecting the actual system configuration (e.g., number of clients, their computing capabilities, the list of currently active patterns) and an **evaluation metrics report** containing the evaluation metrics history of previous rounds (e.g., model accuracy, training time). These files are updated round-by-round to provide a detailed snapshot of the system state and process performance. Leveraging these files, the server initiates the “Architectural Pattern Selection” phase ⑤ by invoking the Multiple AI-Agents layer to determine the updated pattern configuration for the next round.

To augment each AI agent’s decision-making process, a prompt must be carefully designed. A prompt is a textual input given to an LLM that defines its role, specifies the task, provides context if needed, and

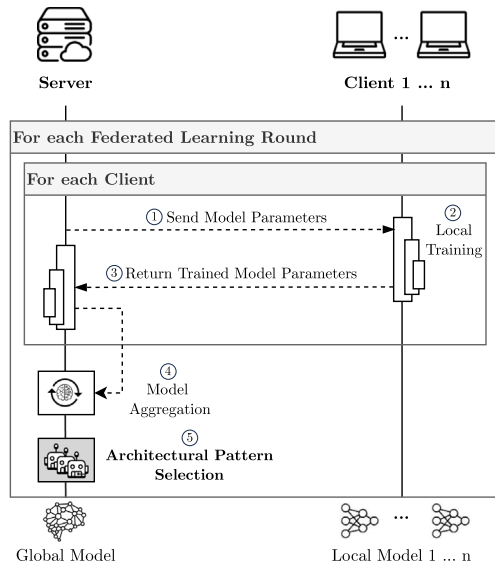


Fig. 7. Federated learning workflow including the multiple AI-agents Layer.

instructs the desired output format (Brown et al., 2020; Liu et al., 2023; Hou et al., 2024) and (Vaswani et al., 2017). Following the guidelines in Bass et al. (2025), we structure the prompt into four parts: (i) *Role* (i.e., the role that the agent must assume), (ii) *Task* (i.e., the task that the agent must perform), (iii) *Guardrails* (i.e., guidelines and rules to avoid failures), (iv) *Output* (i.e., the type of output format expected). To further enhance reasoning capabilities, the agents operate under a Retrieval-Augmented Generation (RAG) configuration. The RAG setup enables LLMs to integrate relevant external knowledge (e.g., database data) at inference time (Arslan et al., 2024). In our implementation, RAG provides access to two key information sources: (i) the system configuration file (see Table 1 for a concrete example), which specifies all the system configuration parameters governing the FL process, and (ii) the evaluation metrics report listing the metrics from all previous rounds. The description of the considered evaluation metrics is reported in Table 2. Both files contain key information that supports the agentic decision to select architectural patterns, ensuring it is driven by runtime contextual parameters and performance metrics.

The final decision is encoded as an updated binary array of ON/OFF flags, which overwrites the corresponding entries in the system configuration file. Once these changes are committed, the server starts a new FL round and propagates the updated configuration to all participants. Consequently, before beginning local training ②, each client applies the received settings by toggling its internal pattern-specific components according to the binary array. This ensures that the subsequent execution of steps ①–④ is performed under a dynamically optimized architectural setup, fully aligned with the agentic decision.

## 5. Experiments design

We aim to address three research questions. We first identify the optimal LLM configurations and prompting strategies to ensure reliable architectural decisions (RQ1). We then evaluate how this orchestration impacts FL system evaluation metrics (RQ2). Finally, we quantify the computational overhead introduced by the Multiple AI-Agents layer to assess the practical sustainability of our approach (RQ3).

The research questions are listed as follows:

**RQ1** Which LLM configurations are most suitable for deploying the Multiple AI-Agents layer?

We evaluate diverse LLM families to identify which model provides the most effective architectural decisions for our use case.

Table 1  
System input parameters.

Parameter	Value
Model Architecture	{CNN, MLP}
Dataset	{FashionMNIST, AG_NEWS}
Data Distribution Type	{IID, non-IID ( $\alpha = 0.5$ )}
Client Data Inflow	{One-shot, Batched}
Network Condition	{Stable, Unstable}
Experiments Configuration Parameters	
(C1) Client 1 - High-Spec.	5 CPU; non-IID; One-shot; Stable
(C2) Client 2 - High-Spec.	5 CPU; non-IID; Batched; Stable
(C3) Client 3 - High-Spec.	5 CPU; IID; Batched; Unstable
(C4) Client 4 - Low-Spec.	3 CPU; non-IID; Batched; Unstable
(C5) Client 5 - Low-Spec.	3 CPU; IID; One-shot; Stable

Table 2  
System evaluation metrics.

Evaluation Metric	Description
MODEL ACCURACY	Global Model Accuracy (F1 Score)
TRAINING TIME	Local Model Training Time per Client
COMMUNICATION TIME	Client-Server Communication Time
TOTAL ROUND TIME	Time to Complete a FL Round
AGENT REASONING TIME	Overhead Introduced by AI-Agents
TOTAL FL TIME	Time to Complete an entire FL Execution

**RQ2** How effective is the Multiple AI-Agents layer in improving FL system evaluation metrics?

We quantify the impact of the Multiple AI-Agents layer on the FL system's evaluation metrics.

**RQ3** What is the computational overhead introduced by the Multiple AI-Agents layer?

We measure the overhead introduced by the Multiple AI-Agents layer on the FL process.

**Subject Systems.** Table 1 reports the input parameters chosen for our experimentation. To evaluate our approach, we consider two different ML tasks for training the global model: image and text classification. For the image classification task, we run 10 FL rounds by training a Convolutional Neural Network (CNN) as the global model on FashionMNIST (Xiao et al., 2017), which contains 70,000 28×28 images (60,000 for training and 10,000 for testing) across 10 clothing categories (e.g., T-shirt/top, trouser, dress, sneaker). For the Text Classification task, we also run 10 FL rounds and train a Multi-Layer Perceptron (MLP) on AG\_NEWS, a benchmark dataset of text news articles labeled by topic (e.g., World, Sport, Business). AG\_NEWS contains 120,000 training samples and 7,600 testing samples. Each sample includes a news title and a short description, which are jointly used as input for text classification. We emulate a heterogeneous federation of five clients with two hardware profiles: *High-Spec.* clients utilize 5 CPU cores, and *Low-Spec.* clients utilize 3 CPU cores. This setup exposes the “Straggler” effect (Vu et al., 2021), where lower-spec clients delay the completion of the round and higher-spec clients are forced to wait idle for them. Client data distribution mixes IID and non-IID using a Dirichlet sampler with  $\alpha = 0.5$ . Data inflow is either *one-shot* (all training data available from round 1) or *batched* (an equal fraction released each round so that by round 10 the client has its full local dataset). Network conditions are *stable* or *unstable*; in the unstable setting, we inject a random client-to-server latency of 20–50 s during model-parameter uploads each round. The concrete per-client setup is: C1 (5 CPU, non-IID, one-shot, stable), C2 (5 CPU, non-IID, batched, stable), C3 (5 CPU, IID, batched, unstable), C4 (3 CPU, non-IID, batched, unstable), and C5 (3 CPU, IID, one-shot, stable). Note that, rather than relying on a synthetic simulation, we use Docker Compose to emulate the FL environment; by containerizing the server and each client with dedicated physical resources, we ensure a more realistic assessment of system evaluation metrics.

**Evaluation Metrics.** Table 2 describes the evaluation metrics used in our experiments, building on the indicators proposed in Jiang et al. (2022) to assess FL systems. MODEL ACCURACY captures the predictive effectiveness of the global model and is reported as F1 score (Wardhani et al., 2019). TRAINING TIME measures the average duration of local model training on the clients, while COMMUNICATION TIME quantifies the time spent exchanging messages between clients and the server. TOTAL ROUND TIME denotes the time required to complete one FL round, and TOTAL FL TIME measures the duration of the entire FL process, i.e., considering all the rounds. The AGENT REASONING TIME measures the additional time required by the Multiple AI-Agents layer to reason over the current FL state and select the architectural patterns for the next round.

**Evaluation Baselines.** To the best of our knowledge, the current literature lacks approaches for the intelligent and dynamic adaptation of FL architectural patterns at runtime. Consequently, to rigorously evaluate the effectiveness of our solution, we define three distinct baseline strategies for comparative analysis. In the (i) Never configuration, all architectural patterns remain disabled throughout the entire training process, serving as a static solution to measure the performance of the FL system in the absence of any architectural pattern. In (ii) Random, each pattern is stochastically activated at each round with a fixed probability of 50%, providing a benchmark to determine whether our approach yields a significant advantage over a non-deterministic adaptation strategy. Finally, (iii) Expert-Driven strategy relies on deterministic activation conditions derived from the quantitative evidence reported in prior empirical studies on FL architectural patterns (Compagnucci et al., 2026b; Lo et al., 2022; Compagnucci et al., 2025). Instead of using fixed absolute thresholds, the activation logic is driven by observable variations in model accuracy, data distribution, and communication overhead. The activation rules at round  $r$  are defined as follows:

$$\begin{aligned} \frac{\text{MODEL ACCURACY}_r}{\text{TOTAL ROUND TIME}_r} < \frac{\text{MODEL ACCURACY}_{r-1}}{\text{TOTAL ROUND TIME}_{r-1}} &\Rightarrow (\text{Client Selector}) \\ \text{JSD}_r > \text{JSD}_{r-1} \wedge \text{MODEL ACCURACY}_r < \text{MODEL ACCURACY}_{r-1} &\Rightarrow (\text{Heterogeneous Data Handler}) \\ \text{COMMUNICATION TIME}_r > \text{COMMUNICATION TIME}_{r-1} &\Rightarrow (\text{Message Compressor}) \end{aligned}$$

The *Client Selector* is activated when the ratio between MODEL ACCURACY and TOTAL ROUND TIME begins to decline; in this case, the pattern discards *Low-Spec* clients to prevent them from creating bottlenecks in the FL process. The *Heterogeneous Data Handler* is triggered when an increase in the JSD metric correlates with a decrease in MODEL ACCURACY, addressing statistical heterogeneity only when it seems to affect model accuracy. Finally, the *MESSAGE COMPRESSOR* is activated when COMMUNICATION TIME increases, specifically when network bottlenecks become an issue on the total round duration.

**Hardware Setup.** Experiments are conducted using a workstation machine with an Intel Xeon W5-2445 chip featuring a 24 Core CPU @3.1 GHz and 64 GB of RAM. Resource allocation across containers is managed via Docker Compose (Docker, Inc., 2026) to emulate a controlled system heterogeneity. For instance, in the *Client Selector* pattern, clients are assigned different CPU quotas to assess performance under varied computational capacities. It is worth noting that the maximum number of containers running at the same time is bounded by the capacity of the host machine (i.e., 24 cores) to avoid CPU overcommitment. This constraint is motivated by the need to avoid exceeding available processing capacity, which may lead to resource contention, affecting execution conditions and compromising the validity of the experiment results (Cohen et al., 2019).

**Statistical Analysis.** To compare the performance of our approach against the baselines, we follow the guidelines proposed by Arcuri and Briand (2011). We repeat each experiment 10 $\times$  and apply the Mann-Whitney U test to assess whether the observed differences are statistically significant. We then compute Vargha and Delaney’s  $A_{12}$  statistic to quantify the effect size of the difference between samples (Vargha and Delaney, 2000). We interpret the effect size using the standard thresholds adopted in the literature: *small* ( $\checkmark$ ) for values

Table 3

LLMs considered in this study.

Model Family	Model Version	Model Parameters	Model Size	Context Length
Llama	3.2	3B	1.3 GB	128K
DeepSeek	r1	8B	5.2 GB	128K
OpenAI	gpt-oss	20B	14 GB	128K

greater than 0.55, *medium* ( $\checkmark\checkmark$ ) for values greater than 0.63, and *large* ( $\checkmark\checkmark\checkmark$ ) for values greater than 0.70. If no statistically significant difference is detected (i.e.,  $0.45 \leq A_{12} \leq 0.55$ ), we use the  $\equiv$  symbol. Note that for metrics where “lower is better” (e.g., execution time), we invert the  $A_{12}$  calculation to ensure that a *Large* effect size consistently represents a significant performance improvement in favor of our approach. Note that the statistical analyses are reported for each RQ and discussed in the corresponding section.

### 5.1. RQ1: Multiple AI-agents layer configuration

This research question aims to explore which configuration settings of AI agents (e.g., LLM type, prompt structure) contribute to optimizing FL systems. We use this experiment to determine which settings yield the optimal configuration. The identified configuration will serve as the LLM baseline model for the subsequent research questions.

We evaluate the agents using three open-source LLMs from the Llama (Llama Team, 2024), DeepSeek (Guo et al., 2025), and OpenAI (Agarwal et al., 2025) families. The choice falls on these model families because, according to the study by Soliman and Keim (2025), they have proven particularly effective in providing accurate and consistent answers to queries involving architectural design decisions. Specifically, the selected models belong to three families: Llama, DeepSeek, and OpenAI. For each family, we consider one representative version, i.e., 3.2 for Llama, r1 for DeepSeek, and gpt-oss for OpenAI. As shown in Table 3, these models also cover different scales in terms of parameter count and storage requirements, ranging from 3B parameters and 1.3 GB for Llama 3.2, to 8B parameters and 5.2 GB for DeepSeek r1, up to 20B parameters and 14 GB for OpenAI gpt-oss. All models share the same 128K-token context length and are evaluated with a temperature of 1.0, following the baseline configuration used in a recent study that investigates the impact of temperature on LLM performance (Li et al., 2025). This supports a consistent setup for comparing different model families and sizes.

We implement and test the LLMs’ reasoning ability by considering two inference strategies that reflect different levels of contextual guidance: (i) the *zero-shot* and (ii) the *few-shot* configurations. In the zero-shot setting, the model relies exclusively on the provided context prompt (Brown et al., 2020; Mosbach et al., 2023). Conversely, the few-shot setting augments the same contextual prompt with a small set of representative examples (Liu et al., 2023). Prompt 1 shows the structure of the few-shot configuration, whereas the zero-shot setting follows the same structure but omits the Examples part.<sup>1</sup> For instance, the contextual prompt includes as guardrails the requirement of having at least two clients participating in each FL round, otherwise the learning is not performed in a federated fashion. As another example, the RAG instructs verification of data from the system configuration file (e.g., CPU/RAM per client contributes to agents’ decisions) and the evaluation metrics report (e.g., the current model accuracy is compared to the previous round to evaluate its variation). Examples include brief explanations of the context and the subsequent decision to be taken. For instance, in Prompt 1 we can see: *Example 1*: “4 High-Spec + 1 Low-Spec” implies the straggler effect (Vu et al., 2021),

<sup>1</sup> All the text prompts used in this work are available at the following link: <https://anonymous.4open.science/r/Agentic-FL-Prompts/README.md>.

**Prompt 1: Simplified Example of the Few-Shot Prompt**

**Contextual Prompt**

**Role:** “You are an expert software architect advising a Federated Learning system [...]”

**Task:** “Recommend which architectural patterns to activate or deactivate for the next round [...]”

**Guardrails:** “At least 2 clients must always participate in the FL round otherwise [...]”

**Output:** “Return only a Rationale text and a JSON object with the following format [...]”

**Retrieval-Augmented Generation (RAG)**

- From the **system configuration file**: extract the actual system parameters (e.g., dataset, CPU/RAM per client, data distr. type)
- From the **evaluation metrics report**: extract the actual system performance (e.g., model accuracy, total round time, training time).

**Examples**

{Example 1: “In an FL system with 4 High-Spec. and 1 Low-Spec. clients, the Low-Spec. client slows down the round, leaving others idle while waiting for synchronization [...]” → Example 1 Decision: “Enable Client Selector (CS=ON) and exclude the Low-Spec. client [...]”}

{Example 2: “CPUs [5, 5, 5, 5, 1]; CS=OFF; high Total Round Time;” → Example 2 Decision: “CS=ON with selection\_value>1; 9× lower Total Round Time;”}

{Example n: ...} → {Example n Decision: ...}

**LLM Output**

**Decision:**{“CS=OFF; HDH=OFF; MC=ON; Rationale=Given the system configuration and the evaluation metrics, I decided to activate the message compressor pattern only because [...]”}

and the corresponding decision is the following: “CS=ON; exclude Low-Spec. client”). It is worth discussing that this scenario can be supported by providing another example that numerically quantifies such system behavior. For instance, Prompt 1 reports, under *Example 2*, the numerical values of CPUs: four clients have 5 CPUs, whereas one client has 1 CPU. The decision is to activate the Client Selector and exclude the client with 1 CPU, resulting in a 9× reduction in Total Round Time. The outcome of the example in Prompt 1 is that LLM decides to select the Message Compressor pattern only.

**Results.** Table 4 shows the results of the LLMs benchmark against baseline configurations. We assess the agent’s performance based on the global MODEL ACCURACY, the total time required to complete the FL process (i.e., TOTAL FL TIME), and the additional overhead introduced by the LLM’s reasoning, in order to identify an effective trade-off between computational efficiency and predictive accuracy. Regarding the baseline results, Never and Random both yield an accuracy of 0.74, but show a large difference in execution time: while Never completes the process in 1666 s, Random has the highest TOTAL FL TIME (3796 seconds), with very high variance. The Expert-Driven baseline, as expected,

shows a slight improvement, reaching an accuracy of  $0.76 \pm 0.01$  and a shorter TOTAL FL TIME of 1619 s.

On the other hand, LLM configurations surpass the baselines in model accuracy, but they do not always achieve shorter total FL times. While agentic approaches are significantly faster than the random baseline (Random), they exhibit a higher TOTAL FL TIME across nearly all configurations compared to static baselines (Never and Expert-Driven), except for llama3.2 models, which achieve a lower TOTAL FL TIME. Zero-Shot (ZS) approaches bring moderate improvements, whereas Few-Shot (FS) settings generally achieve higher accuracy; however, their impact on TOTAL FL TIME depends on the model family. Among them, deepseek-r1 (FS) achieves the highest accuracy, reaching 0.84, with an average total FL process time of 1893 s and an agent reasoning overhead of 18 s. By contrast, llama3.2 (FS) achieves lowest accuracy (i.e., 0.82), completing the entire FL process in an average of 1372 s, with only a minimal reasoning overhead (e.g., 5 s), representing the fastest configuration but with fewer accuracy gains. Finally, gpt-oss (FS) achieves 0.83 accuracy with 1765 s TOTAL FL TIME and 180 s overhead, indicating very high reasoning cost despite a good model accuracy. Statistical analyses of the collected results confirm that all LLM-based configurations significantly outperform the baselines in terms of MODEL ACCURACY, with consistently large effect sizes. For TOTAL FL TIME, only llama3.2 (ZS) and (FS) outperforms the Never and Expert-Driven baselines, whereas the other LLM configurations improve over Random but are slower than Never and Expert-Driven.

In light of these results, we confirm the findings of Soliman and Keim (2025) and select deepseek-r1 with the Few-Shot approach as the baseline LLM configuration for the Multiple AI-Agents layer. Although llama3.2 (FS) is faster, we prioritize the higher accuracy of deepseek-r1 (FS), as even small gains in accuracy are relevant in FL processes (Chen Zhang et al., 2021; Kairouz et al., 2021; Li et al., 2020).

## 5.2. RQ2: Multiple AI-agents layer effectiveness

In this research question, we evaluate the effectiveness of the Multiple AI-Agents layer in optimizing the FL system’s evaluation metrics compared to non-adaptive configurations. After fixing the LLM configuration identified in RQ1, we compare the three coordination strategies (Voting-based, Role-based, and Debate-based) against three baselines (Never, Random, and Expert-Driven). The analysis examines, for both ML tasks, the round-by-round evolution of MODEL ACCURACY, TRAINING TIME, COMMUNICATION TIME, and TOTAL ROUND TIME, thereby assessing both predictive performance and execution efficiency.

**Results.** The performance outcomes are depicted in Fig. 8. Note that all time-related metrics show an upward trend across rounds. This trend is driven by the batched data-inflow regime adopted by some clients: as additional data becomes available at each round, the local dataset grows progressively. Consequently, clients require more time for training, leading to a natural increase in both training and total round time.

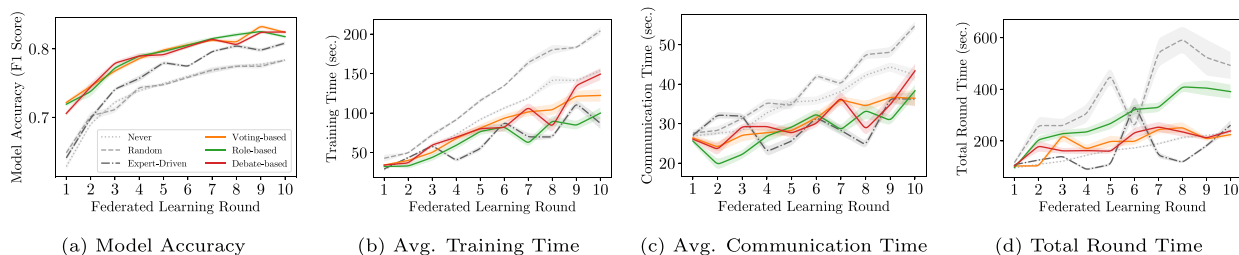
Fig. 8(a) shows the global model accuracy for each configuration in image classification. From the early rounds, the coordination strategies reach higher accuracy and faster convergence than the baselines. The Expert-Driven configuration is the strongest baseline, achieving an F1 score of about 0.79, whereas the agentic strategies reach final values of 0.84–0.85. The three coordination strategies follow similar trends, suggesting that the Multiple AI-Agents layer consistently improves predictive performance in this task. Fig. 8(b) reports the average client training time. The agentic configurations reduce training time with respect to Never and Random. Role-based shows the clearest reduction from round 6 onward, which is consistent with the frequent activation of the Client Selector pattern. However, Expert-Driven remains competitive, maintaining a relatively short training time. Fig. 8(c) reports the average communication time. The agentic strategies reduce communication time compared to Random, whose trend increases

**Table 4**

Never, Random, Expert-Driven, and LLM-based configurations benchmark, together with the statistical analysis for each LLM configuration against the baselines on MODEL ACCURACY and TOTAL FL TIME. (ZS) and (FS) indicate Zero-Shot and Few-Shot prompting strategies adopted by the LLM, respectively.

Configurations	Evaluation metrics			$A_{12}$ on MODEL ACCURACY			$A_{12}$ on TOTAL FL TIME		
	MODEL ACCURACY	TOTAL FL TIME	AGENT OVERHEAD	Never	Random	Expert-Driven	Never	Random	Expert-Driven
Never	0.74 ± 0.00	1666 ± 124	-	-	-	-	-	-	-
Random	0.74 ± 0.01	3796 ± 977	-	-	-	-	-	-	-
Expert-Driven	0.76 ± 0.01	1619 ± 291	-	-	-	-	-	-	-
llama3.2 (ZS)	0.77 ± 0.01	1422 ± 66	7 ± 0	✓✓✓	✓✓✓	✓✓	✓✓✓	✓✓✓	✓✓✓
deepseek-r1 (ZS)	0.80 ± 0.00	1870 ± 240	19 ± 4	✓✓✓	✓✓✓	✓✓✓	XXX	✓✓✓	XXX
gpt-oss (ZS)	0.79 ± 0.01	2276 ± 665	176 ± 6	✓✓✓	✓✓✓	✓✓✓	XXX	✓✓✓	XXX
llama3.2 (FS)	0.82 ± 0.02	1372 ± 151	5 ± 1	✓✓✓	✓✓✓	✓✓✓	✓✓✓	✓✓✓	✓✓✓
deepseek-r1 (FS)	0.84 ± 0.01	1893 ± 188	18 ± 1	✓✓✓	✓✓✓	✓✓✓	XXX	✓✓✓	XXX
gpt-oss (FS)	0.83 ± 0.01	1765 ± 627	180 ± 0	✓✓✓	✓✓✓	✓✓✓	XXX	✓✓✓	XXX

### Image Classification (CNN and FashionMNIST)



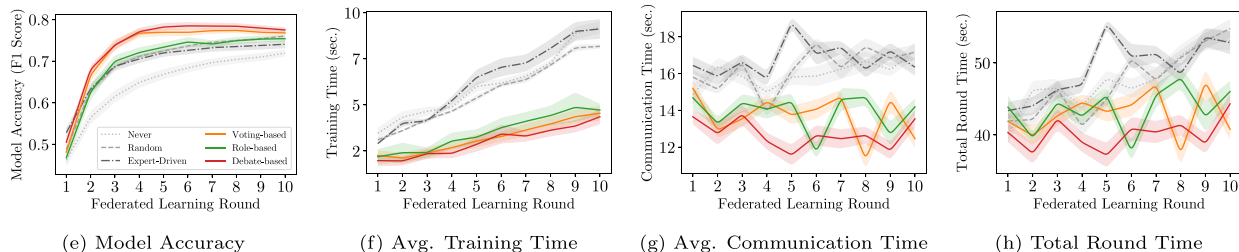
(a) Model Accuracy

(b) Avg. Training Time

(c) Avg. Communication Time

(d) Total Round Time

### Text Classification (MLP and AG NEWS)



(e) Model Accuracy

(f) Avg. Training Time

(g) Avg. Communication Time

(h) Total Round Time

**Fig. 8.** Performance analysis of the Multiple AI-Agents coordination strategies.

across rounds. Voting-based and Role-based remain close to Expert-Driven, while Debate-based is closer to Never and achieves among the lowest communication times. Fig. 8(d) reports the total round time. All agentic strategies improve over Random, but they are not always faster than Never and Expert-Driven. This is expected because Never avoids architectural-pattern overhead, while Expert-Driven activates patterns only under predefined conditions. Among the agentic strategies, Role-based shows the highest total round time, whereas Voting-based and Debate-based follow more similar trends.

For text classification, Fig. 8(e) shows that the differences in model accuracy are less evident. All configurations converge to a similar range, with Debate-based showing a slight advantage. The main benefit of the Multiple AI-Agents layer appears in the time-related metrics. As shown in Figs. 8(f) and 8(g), the coordination strategies reduce both training and communication time with respect to the baselines. For total round time, Fig. 8(h) shows that Role-based and Debate-based improve over all baselines, while Voting-based remains comparable to Never and improves over Random and Expert-Driven.

Table 5 reports the statistical analysis. For image classification, Voting-based significantly improves MODEL ACCURACY over all baselines, while Role-based and Debate-based improve over Never and Random but are statistically equivalent to Expert-Driven. For TRAINING TIME and COMMUNICATION TIME, the agentic strategies generally improve over Never and Random, although Expert-Driven remains competitive. For TOTAL ROUND TIME, all coordination strategies improve over Random, but Expert-Driven remains stronger. For text classification,

accuracy differences are limited, whereas all agentic strategies significantly improve TRAINING TIME and COMMUNICATION TIME. Role-based and Debate-based also improve TOTAL ROUND TIME over all baselines. Overall, the results show that coordinated runtime pattern activation is more effective than random activation.

#### 5.3. RQ3: Multiple AI-agents layer overhead

This research question accounts for the overhead introduced by the Multiple AI-Agents layer. We measure the time required for the agents to reason and elaborate on which pattern or combination should be activated, along with the coordinator (if present), which takes the final decision. Importantly, the AGENT REASONING TIME is not included in the TOTAL FL TIME metric and analyzed in RQ2. We report this overhead for both ML tasks considered in this study: image and text classification. We also perform a statistical analysis to assess whether the observed differences in agent reasoning time among the coordination mechanisms are statistically significant.

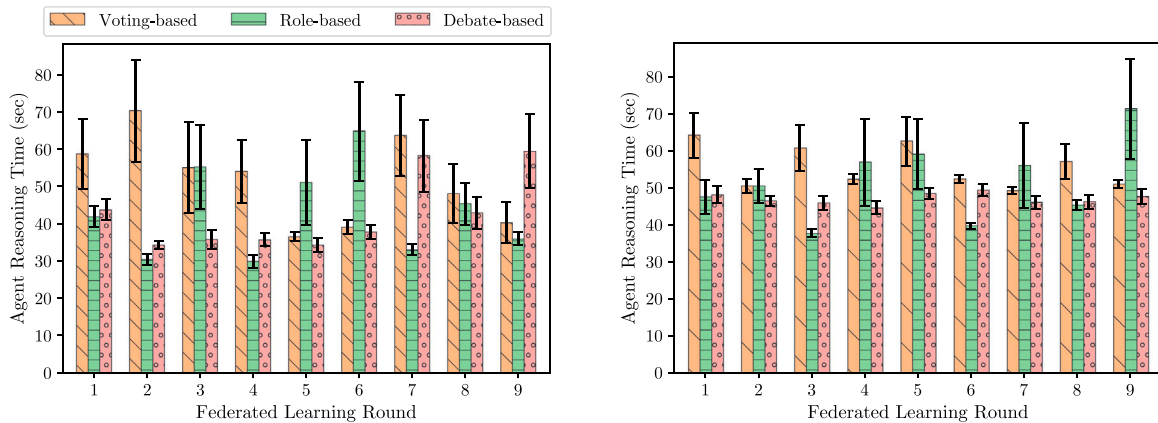
**Results.** Fig. 9 depicts the time overhead introduced by the Multiple AI-Agents layer for both ML tasks in rounds 1–9, since the FL process concludes at the 10th round. For each round, the overhead is computed as the average reasoning time over the 10 repetitions executed with the same configuration and coordination mechanism. Overall, across both ML tasks, the agents' reasoning time per FL round ranges from 30 to 70 s, showing noticeable variability across coordination strategies. On average, for the image recognition task, the Voting-based strategy

**Table 5**  
Statistical analysis for RQ2: Agentic configurations against baselines.

(a) Image Classification (CNN and FashionMNIST)				
Evaluation Metric	Configuration	Never	Random	Expert-Driven
MODEL ACCURACY	Voting-based	✓✓✓	✓✓✓	✓✓✓
	Role-based	✓✓✓	✓✓✓	≡
	Debate-based	✓✓✓	✓✓✓	≡
TRAINING TIME	Voting-based	✓✓✓	✓✓✓	XXX
	Role-based	✓✓✓	✓✓✓	≡
	Debate-based	✓✓✓	✓✓✓	XXX
COMMUNICATION TIME	Voting-based	✓✓✓	✓✓✓	≡
	Role-based	✓✓✓	✓✓✓	≡
	Debate-based	✓✓✓	✓✓✓	≡
TOTAL ROUND TIME	Voting-based	XXX	✓✓✓	XXX
	Role-based	XXX	✓✓✓	XXX
	Debate-based	≡	✓✓✓	XXX

(b) Text Classification (MLP and AG NEWS)				
Evaluation Metric	Configuration	Never	Random	Expert-Driven
MODEL ACCURACY	Voting-based	≡	≡	≡
	Role-based	≡	≡	≡
	Debate-based	✓✓✓	≡	≡
TRAINING TIME	Voting-based	✓✓✓	✓✓✓	✓✓✓
	Role-based	✓✓✓	✓✓✓	✓✓✓
	Debate-based	✓✓✓	✓✓✓	✓✓✓
COMMUNICATION TIME	Voting-based	✓✓✓	✓✓✓	✓✓✓
	Role-based	✓✓✓	✓✓✓	✓✓✓
	Debate-based	✓✓✓	✓✓✓	✓✓✓
TOTAL ROUND TIME	Voting-based	≡	✓✓✓	✓✓✓
	Role-based	✓✓✓	✓✓✓	✓✓✓
	Debate-based	✓✓✓	✓✓✓	✓✓✓



(a) Image Classification (CNN and FashionMNIST).

(b) Text Classification (MLP and AG\_NEWS).

**Fig. 9.** Multiple AI-agents Layer: Average AGENT REASONING TIME Overhead.

**Table 6**  
Statistical analysis for the Multiple AI-Agents Layer REASONING TIME Overhead.

(a) CNN and FashionMNIST			
	Voting-based	Role-based	Debate-based
Voting-based	-	≡	XXX
Role-based	≡	-	≡
Debate-based	XXX	≡	-

(b) MLP and AG_NEWS			
	Voting-based	Role-based	Debate-based
Voting-based	-	≡	XXX
Role-based	≡	-	≡
Debate-based	XXX	≡	-

requires about 51.7 s per round (approximately 465.7 s over 9 rounds), the Role-based strategy about 43.1 s per round (approximately 387.6 s in total), and the Debate-based strategy about 42.2 s per round (approximately 372.0 s over 9 rounds). For text classification, the Voting-based strategy requires about 55.5 s per round (approximately 499.6 s over 9 rounds), the Role-based strategy about 51.5 s per round (approximately 463.7 s in total), and the Debate-based strategy about 46.9 s per round (approximately 422.3 s over 9 rounds). Overall, our experimentation points out that the Debate-based strategy is slightly faster than the others, whereas the Voting-based strategy is the slowest.

Table 6 reports the pairwise statistical analysis of the AGENT REASONING TIME for both ML tasks. The results show a consistent trend across image and text classification tasks. In both cases, Voting-based and Role-based coordination do not show a statistically significant difference. Similarly, Role-based and Debate-based coordination do not show a statistically significant difference. Interestingly, Voting-based coordination is significantly worse than Debate-based coordination with a large effect size in both tasks. This result is consistent with the descriptive results reported above: Voting-based coordination shows the highest average reasoning time, whereas Debate-based coordination shows the lowest.

## 6. Discussion

This section discusses the main findings of our study and reflects on the overall effectiveness of our approach. We analyze how the different coordination strategies influence the FL evaluation metrics and what these results imply for adopting dynamic, agent-driven architectural decisions in practice.

### 6.1. Configuration efficiency score

We evaluate and compare the coordination strategies against the considered baselines using a weighted composite score, denoted by  $S$ . This metric is computed as:

$$S = w_1 \cdot \frac{\text{MODEL ACCURACY} - A_{\min}}{A_{\max} - A_{\min}} + w_2 \cdot \frac{T_{\max} - \text{TOTAL FL TIME}}{T_{\max} - T_{\min}} \quad (1)$$

where  $w_1, w_2 \in [0, 1]$  and  $w_1 + w_2 = 1$ . The weight  $w_1$  determines the contribution of MODEL ACCURACY, whereas  $w_2$  determines the complementary contribution of TOTAL FL TIME. Both terms are normalized in the range  $[0, 1]$ . The second term is inverted so that lower TOTAL FL TIME yields a higher score. Therefore, larger values of  $S$  indicate a more favorable trade-off between predictive quality and execution time. For example, setting  $w_1 = 0.7$  and  $w_2 = 0.3$  means that 70% of the final score is determined by MODEL ACCURACY, whereas the remaining 30% is determined by TOTAL FL TIME. In our evaluation, we adopt the balanced setting with  $w_1 = 0.5$  and  $w_2 = 0.5$ .

The values assumed by this ratio are shown in Fig. 10. Contrary to our initial expectation, not all coordination strategies outperform every baseline. In the image classification task (Fig. 10(a)), Random obtains the lowest efficiency score, indicating that purely stochastic activation of architectural patterns does not yield a favorable trade-off between accuracy and execution time. In the text classification task (Fig. 10(b)), Random remains among the weakest configurations, although the differences among strategies are less marked. In contrast, both Never and Expert-Driven achieve competitive efficiency values in some cases. In particular, Expert-Driven stands out among the baselines, confirming that rule-based activation grounded in empirical evidence can still serve as a strong reference point. At the same time, Never benefits from architectural simplicity, avoiding the additional overhead introduced by pattern management and reducing execution time. Among the agentic approaches, the Role-based strategy yields a less favorable trade-off in the image classification task, while remaining closer to the other agentic configurations in the text classification task. This behavior is mainly explained by its coordination structure: each agent acts as a specialist for a single architectural pattern and

recommends its activation independently from the others. As a result, Role-based tends to activate individual patterns more often, which can increase the Total FL Time and reduce the configuration efficiency score when the additional execution cost is not compensated by proportional accuracy gains. In contrast, Voting-based and Debate-based evaluate pattern configurations more collectively, leading to more competitive median values of  $S$  across the two tasks. Voting-based is particularly effective in image classification (Fig. 10(a)), whereas Debate-based achieves the highest median in text classification (Fig. 10(b)).

Although the configuration efficiency score  $S$  does not imply the existence of a *silver bullet* strategy, the results show that agentic approaches generally achieve higher efficiency, except for the Role-based approach. This confirms that treating architectural pattern activation as a runtime, systematically coordinated decision process can significantly improve the balance between predictive performance and total FL time. At the same time, the magnitude of these improvements depends on the specific system settings (e.g., client heterogeneity, data distribution, and network conditions). In other words, while the agentic approach proves advantageous in our experimental scenarios, its effectiveness may vary across operational conditions and should therefore not be considered always superior in all possible deployment contexts.

### 6.2. Architectural implications: Baseline strategies

The Never configuration represents the simplest architectural setup, where no patterns are activated throughout the FL process. Its main advantage lies in its stability: by avoiding additional processing related to pattern activation, it achieves consistently low total round times. However, this simplicity comes at the cost of limited adaptability, preventing the system from reacting to runtime heterogeneity in data distribution, client resources, or network conditions.

The Random configuration highlights the risks of uncoordinated adaptation. While patterns are dynamically activated, the absence of a decision rationale leads to unstable behavior and the lowest efficiency score. This suggests that architectural variability alone is insufficient; without systematic coordination, dynamic activation may introduce overhead (e.g., by activating architectural patterns when unnecessary) without delivering proportional performance gains.

The Expert-Driven configuration provides a more elaborate alternative. By grounding activation rules in empirical evidence (Compagnucci et al., 2026b, 2025), it achieves competitive efficiency scores and a balanced trade-off between accuracy and total execution time. From an architectural perspective, this confirms that rule-based logic derived from prior studies remains a solid baseline. However, fixed activation rules based only on the metrics from the previous round may limit responsiveness to broader runtime dynamics, especially when decisions would benefit from observing trends over multiple rounds or combining several contextual factors.

### 6.3. Architectural implications: Agentic strategies

The Voting-based strategy achieves a competitive median score  $S$ , especially in image classification, although the Debate-based strategy achieves a higher median in text classification. It also reduces total FL time compared to the Random configuration and avoids the higher activation cost observed in the Role-based configuration. In our setting, each agent reasons about the list of architectural patterns and their impact on the FL metrics, and the voting step aggregates all proposals into a single decision. This behavior aligns with practitioners' expectations for Voting-based cooperation, where decisions made by multiple agents are more accurate and reliable than those made by a single agent (Liu et al., 2025). This strategy leverages *Collective Intelligence*, which filters out locally optimal but globally harmful suggestions and favors pattern combinations that multiple agents independently regard as beneficial, rather than optimizing a single pattern in isolation (Liu et al., 2025).

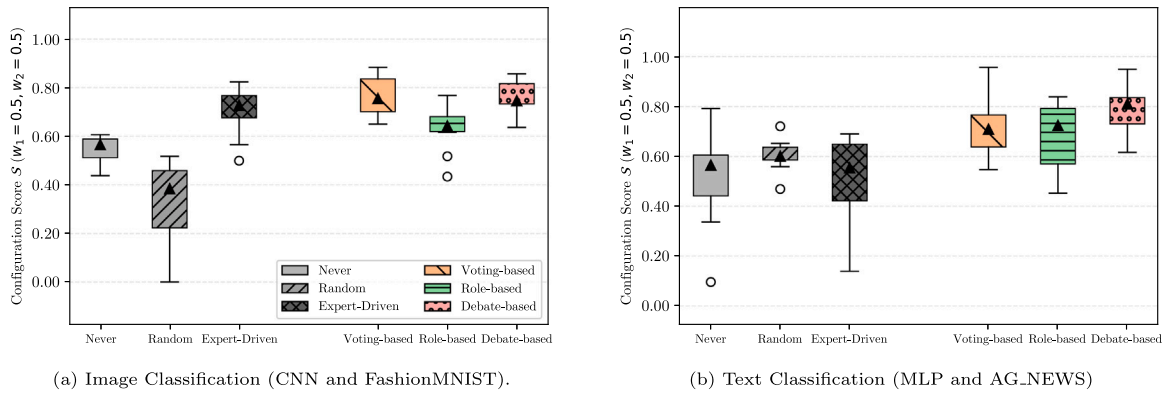


Fig. 10. Comparative analysis of the Configuration Efficiency Score ( $S$ ): Agentic strategies versus baseline configurations.

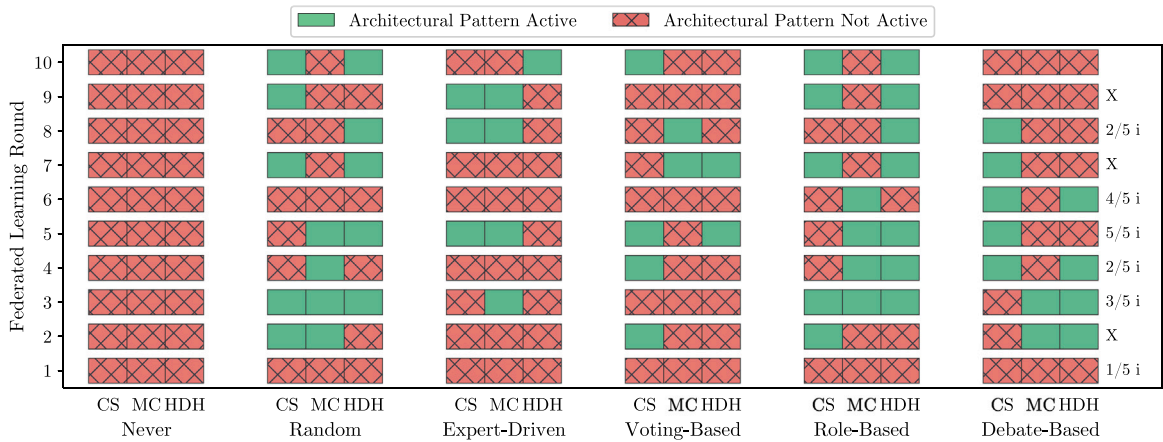


Fig. 11. Architectural Patterns Activation by different strategies. In the Debate-based strategy, labels  $n/5$  indicate the number of debate iterations required to reach consensus, while X marks rounds without consensus, in which the previous configuration of architectural pattern is reused.

The **Role-based** strategy achieves a lower score  $S$  than the Voting- and Debate-based strategies in the image classification task, while remaining closer to the other agentic configurations in the text classification task. It also exhibits higher total FL round time, often close to the Random configuration and clearly higher than the Voting and Debate-based approaches, especially in the image classification task. Accuracy, however, remains broadly comparable to the other strategies; what deteriorates primarily is the total round time, which increases due to frequent pattern activations. This is particularly evident in Fig. 11, where the Role-based strategy shows the highest number of individual FL pattern activations. This behavior stems from each role agent acting as a *Pattern Architect* for a single architectural pattern and reasoning primarily about its own activation, without a full view of how combinations of patterns affect the global FL metrics. Prior work on FL performance analysis (Compagnucci et al., 2025) has shown that combining patterns can yield both benefits and degradations, depending on how they interact. When two or more patterns are activated simultaneously, they may interfere, resulting in degraded performance (e.g., longer total round time). This is why, in this context, software architects are encouraged to use a global coordination strategy that leverages collective intelligence across agents.

The **Debate-based** strategy delivers one of the highest scores  $S$  among the agentic coordination strategies, reaching values comparable to Voting-based in the image classification task and higher median values in the text classification task. This indicates that the explicit exchange of arguments among agents can help identify effective pattern combinations, especially when the decision involves multiple interacting architectural concerns. However, the iterative nature of the debate still introduces variability: when agents disagree for several turns, the

system spends more time debating before making a decision, thereby directly increasing the total FL time. This behavior is consistent with practitioners’ reports on Debate-based strategy (Liu et al., 2025), where richer debates among agents introduce extra communication and coordination time. In our setting, Debate-based is therefore a competitive coordination strategy rather than only an occasional best-case solution: it can occasionally outperform all other strategies when consensus is reached quickly; however, it can also take a long time to reach a decision when consensus is not reached quickly. Software architects should consider the Debate-based strategy for scenarios where higher decision quality is required and coordination variability is acceptable. Future work will investigate the Debate-based strategy with more capable LLMs (e.g., models with more parameters) to assess whether higher-capacity agents can shorten the time to reach consensus and, in turn, reduce variability, making this strategy a more reliable option.

#### 6.4. Adaptation overhead

The time overhead introduced by the Multiple AI-Agents layer is analyzed by measuring the reasoning cost required to decide which architectural patterns to activate or deactivate at runtime. It affects the FL execution in two different ways: (i) it can reduce the time required to complete the overall FL execution by avoiding unnecessary pattern activations, but (ii) it also introduces an additional cost due to LLM reasoning.

From the **TOTAL FL TIME** perspective, the agentic strategies are effective in reducing the time required to complete the overall FL execution by making pattern activation more selective, especially compared with

the Random baseline. In the image classification task, Random requires, on average, about 379.1 s per round, whereas Voting-based and Debate-based coordination reduce this value to about 192.2 and 193.9 s, corresponding to reductions of approximately 49.3% and 48.8%, respectively. Role-based coordination also improves over Random, with an average round time of about 288.5 s, corresponding to a reduction of approximately 23.9%. A similar effect is observed in the text classification task: Voting-based, Role-based, and Debate-based coordination require about 42.9, 43.6, and 40.2 s per round, respectively, compared to 47.3 s for Random and 49.3 s for Expert-Driven.

However, the reduction in TOTAL FL TIME comes with an additional cost. The LLM reasoning phase introduces extra execution time, since agents need to analyze the current FL state and decide the architectural configuration for the next round. Across the two ML tasks, the average reasoning overhead is about 53.6 s per round for Voting-based coordination, 47.3 s per round for Role-based coordination, and 44.6 s per round for Debate-based coordination. Therefore, although agentic configurations reduce the TOTAL FL TIME through more informed architectural decisions, accounting for the agentic adaptation phase increases overall execution time, adding an average cost of about 33% over the TOTAL FL TIME. However, this overhead should be interpreted together with the effectiveness gains of the Multiple AI-Agents layer. Across all experiments, agentic configurations improve average model accuracy relative to the baselines, indicating that the reasoning cost yields measurable predictive gains.

Overall, these results clarify the main trade-off introduced by the Multiple AI-Agents layer. On the one hand, the agents improve the quality of architectural decisions: they reduce unnecessary time overhead caused by architectural patterns and avoid ineffective pattern combinations. On the other hand, they introduce a new source of overhead due to LLM reasoning. Therefore, the Multiple AI-Agents layer shifts part of the execution cost from pattern execution to decision making. This cost can be justified when reasoning prevents the activation of expensive patterns that would not improve the FL process (e.g., avoiding the Heterogeneous Data Handler when synthetic data generation would increase training time without improving accuracy). In light of this, reasoning overhead should not be interpreted in isolation. The absence of large differences in reasoning time among coordination strategies suggests that this metric alone is insufficient to determine the most suitable strategy. Instead, the choice of the coordination mechanism should be considered together with the effectiveness results discussed in RQ2.

### 6.5. Threats to validity

**External Validity.** Generalization of findings is not guaranteed since we present a specific experimental setup. Our current evaluation is limited to a restricted setup of the considered architectural patterns, the number of datasets, FL rounds, and clients. This setup reflects a methodological trade-off, since we prioritize: (i) the implementation of patterns contributing to different evaluation metrics; (ii) the datasets learning different data types; (iii) the hardware-level fidelity by assigning dedicated CPU resources to each client container and avoiding CPU overcommitment. Exceeding available processing capacity can lead to resource contention, potentially affecting execution conditions and compromising the correctness of the collected metrics (Cohen et al., 2019). While this limits the number of concurrent clients, it endorses the reliability of the performance metrics used by the agents for runtime architectural decisions. Yet, there is a combination of diverse LLMs, coordination strategies, and FL scenarios spanning multiple dimensions, including LLM families and prompting regimes (zero-shot and few-shot), datasets performing different classification tasks, heterogeneous client hardware profiles, and varying network conditions. Such variations span a broad range of variegated conditions and provide initial evidence that our approach can be applied beyond a single case. To

further strengthen external validity, we plan to expand the experimentation by considering more architectural patterns, as well as further datasets and models, along with more FL rounds and clients.

**Internal Validity.** The settings and input parameters adopted for performance analysis can be vulnerable to potential threats, as they rely on numerical values that serve as evidence of performance variations. Yet, identifying appropriate configuration settings and parameter values remains a challenge in software performance (Bondi, 2015), and additional alternatives could be explored in future work. To this end, we make the framework available to replicate the presented experimental results and potentially investigate additional variations of the input parameters. Besides, adopting computational power as a client selection criterion may represent a threat, since excluding low-power devices can improve overall system efficiency at the cost of unfairly excluding certain types of clients that may contain relevant training data (Lo et al., 2022), possibly adding bias to the overall model. As future work, we plan to develop different client selection strategies, e.g., by measuring the fairness of client selection during training without compromising their privacy.

**Construct Validity.** Threats to construct validity concern possible misinterpretation of the measured metrics. We assess FL systems using standard metrics, e.g., F1 score. We do not employ text generation metrics such as ROUGE-1 (Lin, 2004) or BERTScore (Tianyi Zhang et al., 2020), as our setting lacks ground-truth outputs for comparison with model responses. Instead, we evaluate LLMs indirectly through FL metrics as proxies for their decision effectiveness. Moreover, a significant threat in this context is the inherent non-determinism of LLMs, which can introduce variability in the agents' decisions across different runs, potentially masking the actual impact of the coordination strategies. To mitigate this, we repeat each experiment 10x and assess the statistical significance of the observed variations using the Mann-Whitney U test and Vargha and Delaney's  $\hat{A}_{12}$  effect size.

## 7. Conclusion

This paper investigates how architectural decisions in FL systems can be elevated from static design-time choices to dynamic runtime mechanisms. We introduce an agentic framework in which multiple AI agents, coordinated through Voting-, Role-, and Debate-based strategies, dynamically (de)activate three architectural patterns: Client Selector, Heterogeneous Data Handler, and Message Compressor. Our results show that treating architectural pattern activation as a systematically coordinated runtime decision can improve the balance between predictive performance and total execution time. In particular, agentic coordination strategies achieve higher configuration efficiency scores in almost all cases in our experimental setting, showing that context-aware architectural reconfiguration improves the trade-off between predictive performance and core FL execution time, while the additional reasoning overhead remains a practical cost to be considered. These results highlight a complementary optimization dimension for FL systems: beyond tuning training parameters or global model aggregation strategies, the system architecture itself can be dynamically reconfigured based on evolving runtime conditions and historical performance metrics. While no single coordination strategy is a *silver bullet* across all scenarios, our findings provide quantitative evidence that agent-driven architectural adaptation can significantly improve the efficiency of FL systems under heterogeneous, dynamic operating conditions.

As future work, we plan to investigate additional agentic coordination strategies and experiment with new families of LLMs (e.g., Gemma 4 Google DeepMind, 2026, Mistral-3 Mistral AI, 2026, and Qwen Qwen Team, 2025) to further assess their reasoning capabilities. This is particularly relevant, as LLMs are advancing rapidly, with new models and updated versions released frequently. Moreover, we intend to integrate a broader set of FL architectural patterns in order to expand the space of runtime design alternatives.

## CRedit authorship contribution statement

**Ivan Compagnucci:** Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Qinghua Lu:** Writing – review & editing, Visualization, Validation, Supervision, Conceptualization. **Catia Trubiani:** Writing – review & editing, Writing – original draft, Visualization, Validation, Supervision, Project administration, Methodology, Investigation, Funding acquisition, Formal analysis, Conceptualization.

## Declaration of Generative AI use

During the preparation of this work, the authors used GPT-5.5 and Grammarly for grammar and spelling checks. After using these tools, the authors reviewed and edited the content as needed and took full responsibility for the publication's content. LLMs (i.e., llama3.2, deepseek-r1, and gpt-oss) are integrated into the Multiple AI-Agent layer as a core component of the proposed framework and empirically evaluated.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgments

This work has been partially funded by the Italian Ministry of University and Research (MUR) Department of Excellence 2023–2027 for GSSI, and the MUR-PRIN project 20228FT78M DREAM (modular software Design to Reduce uncertainty in Ethics-based cyber-physical systems).

## Data availability

The results of the experiments are available on Zenodo [Compagnucci et al. \(2026a\)](#).

## References

- Agarwal, Sandhini, et al., 2025. gpt-oss-120b & gpt-oss-20b model card. CoRR abs/2508.10925.
- Amershi, Saleema, Begel, Andrew, Bird, Christian, DeLine, Robert, Gall, Harald, Kamar, Ece, Nagappan, Nachiappan, Nushi, Besmira, Zimmermann, Thomas, 2019. Software engineering for machine learning: A case study. In: International Conference on Software Engineering: Software Engineering in Practice. ICSE-SEIP, pp. 291–300.
- Arcuri, Andrea, Briand, Lionel, 2011. A practical guide for using statistical tests to assess randomized algorithms in software engineering. In: International Conference on Software Engineering. ICSE, pp. 1–10.
- Arslan, Muhammad, Ghanem, Hussam, Munawar, Saba, Cruz, Christophe, 2024. A survey on RAG with LLMs. In: Knowledge-Based and Intelligent Information & Engineering Systems, vol. 246, pp. 3781–3790.
- Arun, Shrikara, Tedla, Meghana, Vaidhyathan, Karthik, 2025. LLMs for generation of architectural components: An exploratory empirical study in the serverless world. In: International Conference on Software Architecture. ICSA, pp. 25–36.
- Baresi, Luciano, Lestingi, Livia, Wehbe, Iyad, 2025. Architecting federated learning systems: A requirement-driven approach. ECSA, In: European Conference on Software Architecture, vol. 15929, Springer, pp. 224–239.
- Baresi, Luciano, Quattrocchi, Giovanni, Rasi, Nicholas, 2021. Federated machine learning as a self-adaptive problem. In: International Symposium on Software Engineering for Adaptive and Self-Managing Systems. SEAMS, IEEE, pp. 41–47.
- Bass, Len, Lu, Qinghua, Weber, Ingo, Zhu, Liming, 2025. Engineering AI systems: Architecture and DevOps essentials.
- Bondi, André B., 2015. Foundations of software and system performance engineering: Process, performance modeling, requirements, testing, scalability, and practice.
- Briggs, Christopher, Fan, Zhong, Andras, Peter, 2020. Federated learning with hierarchical clustering of local updates to improve training on non-IID data. In: International Conference on Neural Network. IJCNN, pp. 1–9.

- Brown, Tom B., Mann, Benjamin, Ryder, Nick, Subbiah, Melanie, Kaplan, Jared, Dhariwal, Prafulla, Neelakantan, Arvind, Shyam, Pranav, Sastry, Girish, Askell, Amanda, Agarwal, Sandhini, Herbert-Voss, Ariel, Krueger, Gretchen, Henighan, Tom, Child, Rewon, Ramesh, Aditya, Ziegler, Daniel M., Wu, Jeffrey, Winter, Clemens, Hesse, Christopher, Chen, Mark, Sigler, Eric, Litwin, Mateusz, Gray, Scott, Chess, Benjamin, Clark, Jack, Berner, Christopher, McCandlish, Sam, Radford, Alec, Sutskever, Ilya, Amodei, Dario, 2020. Language models are few-shot learners. In: Conference on Neural Information Processing Systems. NeurIPS.
- Brun, Yuriy, Serugendo, Giovanna Di Marzo, Gacek, Cristina, Giese, Holger, Kienle, Holger M., Litoiu, Marin, Müller, Hausi A., Pezzè, Mauro, Shaw, Mary, 2009. Engineering self-adaptive systems through feedback loops. In: Software Engineering for Self-Adaptive Systems. pp. 48–70.
- Cohen, Maxime C., Keller, Philipp W., Mirrokni, Vahab S., Zadimoghaddam, Morteza, 2019. Overcommitment in cloud services: Bin packing with chance constraints. Manag. Sci. 65 (7), 3255–3271.
- Compagnucci, Ivan, Lu, Qinghua, Trubiani, Catia, 2026a. Open science artifact: Agentic runtime reconfiguration of architectural patterns in federated learning. <http://dx.doi.org/10.5281/zenodo.20006358>.
- Compagnucci, Ivan, Pinciroli, Riccardo, Trubiani, Catia, 2025. Performance analysis of architectural patterns for federated learning systems. In: International Conference on Software Architecture. ICSA, pp. 289–300.
- Compagnucci, Ivan, Pinciroli, Riccardo, Trubiani, Catia, 2026b. Experimenting architectural patterns in federated learning systems. J. Syst. Softw. 232, 112655.
- Compagnucci, Ivan, Trubiani, Catia, 2025. Towards AI agents for selecting architectural patterns in federated learning systems. In: Proceedings of the Fourth Conference on System and Service Quality (QualITA 2025), vol. 4080.
- Dayan, Ittai, et al., 2021. Federated learning for predicting clinical outcomes in patients with COVID-19. Nature Med. 27 (10), 1735–1743.
- Dhar, Rudra, Vaidhyathan, Karthik, Varma, Vasudeva, 2024a. Can LLMs generate architectural design decisions? - an exploratory empirical study. In: International Conference on Software Architecture. ICSA, pp. 79–89.
- Dhar, Rudra, Vaidhyathan, Karthik, Varma, Vasudeva, 2024b. Leveraging generative AI for architecture knowledge management. In: International Conference on Software Architecture Companion. ICSA-C, IEEE, pp. 163–166.
- Docker, Inc., 2026. Docker. <https://www.docker.com/>. (Accessed 09 May 2026).
- Fu, Lei, Zhang, Huanle, Gao, Ge, Zhang, Mi, Liu, Xin, 2023. Client selection in federated learning: Principles, challenges, and opportunities. IEEE Internet Things J. 10 (24), 21811–21819.
- Goodfellow, Ian J., Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair, Sherjil, Courville, Aaron C., Bengio, Yoshua, 2014. Generative adversarial nets. In: Conference on Neural Information Processing Systems. NeurIPS, pp. 2672–2680.
- Google DeepMind, 2026. Gemma 4 model card. [https://ai.google.dev/gemma/docs/core/model\\_card\\_4](https://ai.google.dev/gemma/docs/core/model_card_4). (Accessed 09 May 2026).
- Guo, Daya, et al., 2025. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. CoRR abs/2501.12948.
- Hou, Xinyi, Zhao, Yanjie, Liu, Yue, Yang, Zhou, Wang, Kailong, Li, Li, Luo, Xiapu, Lo, David, Grundy, John, Wang, Haoyu, 2024. Large language models for software engineering: A systematic literature review. ACM Trans. Softw. Eng. Methodol. 33 (8), 220:1–220:79.
- Hu, Zhiyao, Li, Dongsheng, Yang, Ke, Xu, Ying, Peng, Baoyun, 2025. Optimizing data distributions based on jensen-Shannon divergence for federated learning. Tsinghua Sci. Technol. 30 (2), 670–681.
- Ilhan, Fatih, Su, Gong, Liu, Ling, 2023. Scalefl: Resource-adaptive federated learning with heterogeneous clients. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 24532–24541.
- Jiang, Yuang, Wang, Shiqiang, Valls, Victor, Ko, Bong Jun, Lee, Wei-Han, Leung, Kin K, Tassiulas, Leandros, 2022. Model pruning enables efficient federated learning on edge devices. IEEE Trans. Neural Networks Learn. Syst. 34 (12), 10374–10386.
- Kairouz, Peter, et al., 2021. Advances and open problems in federated learning. Found. Trends Mach. Learn. 14 (1–2), 1–210.
- Kephart, Jeffrey O., Chess, David M., 2003. The vision of autonomic computing. Comput. 36 (1), 41–50.
- Lai, Fan, Dai, Yinwei, Singapuram, Sanjay, Liu, Jiachen, Zhu, Xiangfeng, Madhyastha, Harsha, Chowdhury, Mosharaf, 2022. FedScale: Benchmarking model and system performance of federated learning at scale. In: Proceedings of Machine Learning Research. PMLR, pp. 11814–11827.
- Li, Li, Duan, Moming, Liu, Duo, Zhang, Yu, Ren, Ao, Chen, Xianzhang, Tan, Yujian, Wang, Chengliang, 2021. FedSAE: A novel self-adaptive federated learning framework in heterogeneous systems. In: Proceedings of the International Joint Conference on Neural Networks. IJCNN, pp. 1–10.
- Li, Tian, Sahu, Anit Kumar, Talwalkar, Ameet, Smith, Virginia, 2020. Federated learning: Challenges, methods, and future directions. IEEE Signal Process. Mag. 37 (3), 50–60.
- Li, Lujun, Sleem, Lama, Gentile, Niccolò, Nichil, Geoffrey, State, Radu, 2025. Exploring the impact of temperature on large language models: Hot or cold? Procedia Comput. Sci. 264, 242–251.
- Li, Jialong, Zhang, Mingyue, Li, Nianyu, Weyns, Danny, Jin, Zhi, Tei, Kenji, 2024. Exploring the potential of large language models in self-adaptive systems. In: Proceedings of the International Symposium on Software Engineering for Adaptive and Self-Managing Systems. SEAMS, pp. 77–83.

- Lin, Chin-Yew, 2004. ROUGE: A package for automatic evaluation of summaries. In: Text Summarization Branches Out. pp. 74–81.
- Liu, Yue, Lo, Sin Kit, Lu, Qinghua, Zhu, Liming, Zhao, Dehai, Xu, Xiwei, Harrer, Stefan, Whittle, Jon, 2025. Agent design pattern catalogue: A collection of architectural patterns for foundation model based agents. *J. Syst. Softw.* 220, 112278.
- Liu, Pengfei, Yuan, Weizhe, Fu, Jinlan, Jiang, Zhengbao, Hayashi, Hiroaki, Neubig, Graham, 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Comput. Surv.* 55 (9), 195:1–195:35.
- Llama Team, 2024. The llama 3 herd of models. *CoRR abs/2407.21783*.
- Lo, Sin Kit, Lu, Qinghua, Paik, Hye-Young, Zhu, Liming, 2021. FLRA: A reference architecture for federated learning systems. In: *European Conference on Software Architecture*. Springer, pp. 83–98.
- Lo, Sin Kit, Lu, Qinghua, Zhu, Liming, Paik, Hye-Young, Xu, Xiwei, Wang, Chen, 2022. Architectural patterns for the design of federated learning systems. *J. Syst. Softw.* 191, 111357.
- Mayhoub, Samara, M. Shami, Tareq, 2024. A review of client selection methods in federated learning. *Arch. Comput. Methods Eng.* 31 (2), 1129–1152.
- McMahan, Brendan, Moore, Eider, Ramage, Daniel, Hampson, Seth, y Arcas, Blaise Agüera, 2017. Communication-efficient learning of deep networks from decentralized data. In: *International Conference on Artificial Intelligence and Statistics*, AISTATS, vol. 54, PMLR, pp. 1273–1282.
- Menon, Sindhu, Addula, Santosh Reddy, Parkavi, A., Subbalakshmi, Ch., Dhandayuthapani, V. Bala, Pokkuluri, Kiran Sree, Soni, Anita, 2024. Streamlining task planning systems for improved enactment in contemporary computing surroundings. *SN Comput. Sci.* 5 (8), 993.
- Mistral AI, 2026. Ministral 3. *CoRR abs/2601.08584*.
- Mosbach, Marius, Pimentel, Tiago, Ravfogel, Shauli, Klakow, Dietrich, Elazar, Yanai, 2023. Few-shot fine-tuning vs. In-context learning: A fair comparison and evaluation. In: *Findings of the Association for Computational Linguistics: ACL*. pp. 12284–12314.
- Pace, Jorge Andrés Díaz, Tommasel, Antonela, Capilla, Rafael, 2024. Helping novice architects to make quality design decisions using an LLM-based assistant. In: *Proceedings of European Conference on Software Architecture (ECSA)*, vol. 14889, pp. 324–332.
- Pace, Jorge Andrés Díaz, Tommasel, Antonela, Capilla, Rafael, Ramírez, Yamid E., 2025. Architecture exploration and reflection meet LLM-based agents. In: *International Conference on Software Architecture Companion*. ICSA-C, IEEE, pp. 1–5.
- Qwen Team, 2025. Qwen3 technical report. *CoRR abs/2505.09388*.
- Richards, Mark, 2015. *Software Architecture Patterns*, vol. 4.
- Rizk, Elsa, Vlaski, Stefan, Sayed, Ali H., 2020. Dynamic federated learning. In: *Proceedings of the International Workshop on Signal Processing Advances in Wireless Communications*. SPAWC, pp. 1–5.
- Sánchez, Pedro Miguel Sánchez, Celdrán, Alberto Huertas, Xie, Ning, Bovet, Jérôme, Pérez, Gregorio Martínez, Stiller, Burkhard, 2024. FederatedTrust: A solution for trustworthy federated learning. *Future Gener. Comput. Syst.* 152, 83–98.
- Singh, Neha, Adhikari, Mainak, 2025. SelfFed: Self-adaptive federated learning with non-IID data on heterogeneous edge devices for bias mitigation and enhance training efficiency. *Inf. Fusion* 118, 102932.
- Soliman, Mohamed, Keim, Jan, 2025. Do large language models contain software architectural knowledge? : An exploratory case study with GPT. In: *International Conference on Software Architecture*. ICSA, IEEE, pp. 13–24.
- Vaidhyanathan, Karthik, Muccini, Henry, 2025. Software architecture in the age of agentic AI. *ECSA*, In: *European Conference on Software Workshops*, vol. 15982, Springer, pp. 41–49.
- Vargha, András, Delaney, Harold D., 2000. A critique and improvement of the CL common language effect size statistics of McGraw and wong. *J. Educ. Behavioral Stat.* 25 (2), 101–132.
- Vaswani, Ashish, Shazeer, Noam, Parmar, Niki, Uszkoreit, Jakob, Jones, Llion, Gomez, Aidan N., Kaiser, Lukasz, Polosukhin, Illia, 2017. Attention is all you need. In: *Conference on Neural Information Processing Systems*. NeurIPS, pp. 5998–6008.
- Verbraeken, Joost, Wolting, Matthijs, Katzy, Jonathan, Kloppenburg, Jeroen, Verbeelen, Tim, Rellermeyer, Jan S., 2021. A survey on distributed machine learning. *ACM Comput. Surv.* 53 (2), 30:1–30:33.
- Vu, Tung T., Ngo, Duy T., Ngo, Hien Quoc, Dao, Minh N., Tran, Nguyen H., Middleton, Richard H., 2021. Straggler effect mitigation for federated learning in cell-free massive MIMO. In: *IEEE International Conference on Communications*. ICC, pp. 1–6.
- Wang, Shiqiang, Tuor, Tiffany, Salonidis, Theodoros, Leung, Kin K., Makaya, Christian, He, Ting, Chan, Kevin, 2019. Adaptive federated learning in resource constrained edge computing systems. *IEEE J. Sel. Areas Commun.* 37 (6), 1205–1221.
- Wardhani, Ni Wayan Surya, Rochayani, Masithoh Yessi, Iriany, Atiek, Sulistyono, Agus Dwi, Lestantyo, Prayudi, 2019. Cross-validation metrics for evaluating classification performance on imbalanced data. In: *International Conference on Computer, Control, Informatics and Its Applications*. IC3INA, pp. 14–18.
- Xiao, Han, Rasul, Kashif, Vollgraf, Roland, 2017. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *CoRR abs/1708.07747*.
- Yang, Qiang, Liu, Yang, Chen, Tianjian, Tong, Yongxin, 2019. Federated machine learning: Concept and applications. *ACM Trans. Intell. Syst. Technol.* 10 (2).
- Yu, Lantao, Zhang, Weinan, Wang, Jun, Yu, Yong, 2017. Seqgan: Sequence generative adversarial nets with policy gradient. In: *Conference on Artificial Intelligence AAAI*. pp. 2852–2858.
- Zhang, Hongyi, Bosch, Jan, Olsson, Helena Holmström, 2020. Federated learning systems: Architecture alternatives. In: *Asia-Pacific Software Engineering Conference*. APSEC, pp. 385–394.
- Zhang, Yizhe, Gan, Zhe, Fan, Kai, Chen, Zhi, Henao, Ricardo, Shen, Dinghan, Carin, Lawrence, 2017. Adversarial feature matching for text generation. In: *International Conference on Machine Learning*, ICML, vol. 70, PMLR, pp. 4006–4015.
- Zhang, Jie, Guo, Song, Qu, Zhihao, Zeng, Deze, Zhan, Yufeng, Liu, Qifeng, Akerkar, Rajendra, 2021. Adaptive federated learning on non-iid data with resource constraint. *IEEE Trans. Comput.* 71 (7), 1655–1667.
- Zhang, Tianyi, Kishore, Varsha, Wu, Felix, Weinberger, Kilian Q., Artzi, Yoav, 2020. BERTScore: Evaluating text generation with BERT. In: *International Conference on Learning Representations*. ICLR.
- Zhang, Chen, Xie, Yu, Bai, Hang, Yu, Bin, Li, Weihong, Gao, Yuan, 2021. A survey on federated learning. *Knowl.-Based Syst.* 216, 106775.