(REVIEW ARTICLE)

# Understanding data heterogeneity in federated learning

Kuldeep Deshwal *

*Proofpoint Inc, USA.*

## Abstract

Federated learning enables machine learning across distributed devices without centralizing sensitive data, preserving privacy while creating intelligent systems from collective knowledge. Data heterogeneity, the natural variation in information across participating devices, presents significant challenges including convergence instability, model bias, communication inefficiency, privacy-utility tradeoffs, and computational imbalance. Despite these obstacles, heterogeneity offers advantages like improved model generalization, personalization opportunities, greater real-world applicability, enhanced privacy protection, and better fault tolerance when properly managed. Current solutions address these challenges through personalized federated learning, robust aggregation methods, federated distillation, client clustering, and adaptive participation strategies, while future directions focus on developing advanced heterogeneity metrics, cross-organizational techniques, dynamic adaptation mechanisms, hardware-aware algorithms, theoretical foundations, and standardized benchmarks to further enhance performance in diverse data environments.

**Keywords:** Adaptation; Decentralization; Heterogeneity; Personalization; Privacy

## 1. Introduction

Think about all the personal information stored on your smartphone - your texts, photos, health data, and browsing habits. Now imagine if we could build smart systems that learn from this valuable information without ever seeing your actual data. This is exactly what Federated Learning does. Instead of the traditional approach where all data is uploaded to a central server for processing, federated learning flips the script by bringing the learning process directly to your device. The central server sends a starting model to your phone and thousands of other devices. Each device then trains this model using only its local data. Once training is complete, only the improvements to the model are sent back to the central server - never your personal information. The server combines these improvements from all participating devices to create a better global model, which is then shared with everyone. This cycle continues, making the system smarter over time while keeping your data safely on your device.

This innovative approach solves major privacy concerns since sensitive information never leaves your phone. It also dramatically reduces the amount of data that needs to be transferred over networks, saving bandwidth and battery life. For example, when your keyboard app suggests text completions, it can learn your writing style directly on your phone without sending your messages to external servers. Similarly, voice assistants can improve their recognition of your specific accent without recording and uploading your voice. Healthcare applications can analyze patient data across multiple hospitals without sharing confidential medical records. The beauty of federated learning is that it allows collaboration without compromising individual privacy. This privacy-preserving approach was pioneered by Mcmahan et al [1], who introduced FederatedAveraging, the first efficient algorithm for training deep networks on decentralized data without sharing raw information.

---

* Corresponding author: Kuldeep Deshwal

However, federated learning faces a significant hurdle called data heterogeneity. This simply means that the information on different devices varies tremendously. Consider how differently people use their smartphones: some primarily text in short bursts during work hours, others mainly share photos and videos on weekends, and some regularly use voice commands in various languages. The data generated by these different usage patterns isn't uniform - it's heterogeneous. This variation extends beyond just user behavior to include differences in device types, geographical locations, and cultural contexts. A model that works perfectly for users in urban America might perform poorly for those in rural India due to these differences. This data heterogeneity makes building effective models much more complicated, as the system must somehow create a global model that works well across incredibly diverse usage patterns and data types without ever seeing the raw data itself.

## 2. Understanding Data Heterogeneity

Data heterogeneity in federated learning is a technical way of saying that the information stored across different devices isn't uniform or consistent. When experts talk about "non-independent and identically distributed (non-IID) data," they're simply pointing out that the data on one person's phone looks very different from the data on someone else's phone. This is quite different from traditional machine learning, where researchers carefully collect balanced datasets that represent all the scenarios a model might encounter. In the real world of federated learning, each device contains a unique slice of information based on how that specific person uses their device, creating natural but challenging differences that the learning system must handle. Researchers formalized this concept by defining non-IID data in federated settings as having statistical heterogeneity where local data distributions differ substantially from the population distribution [2].

### 2.1. What Makes Data Heterogeneous?

To understand data heterogeneity better, think about three friends who all use the same messaging app. Alex works in an office and mainly sends quick text messages during the day, keeping communications brief and professional. These messages might contain work jargon, short replies, and are typically sent Monday through Friday during business hours. Blake is more of a social media enthusiast who rarely texts but frequently shares photos and videos on weekends. Blake's data consists primarily of image and video files with occasional captions, mostly created during evening hours and weekends. Casey, who has family in different countries, often sends voice messages in multiple languages. Casey's data includes long audio files in various languages that might be sent at any time of day to accommodate different time zones.

If all three friends participate in federated learning for a text prediction feature, the system faces a major challenge. Alex's device contains mostly short, professional text with predictable patterns. Blake's device has very little text data but lots of media with captions that might use more casual language. Casey's device contains voice data in multiple languages that would first need to be converted to text. This creates several types of variation. Data distribution varies because the content types differ dramatically between users. Data quantity differs as Alex might send hundreds of messages daily while Blake might only send a few text captions per week. Data quality varies because Alex's text messages are clear and complete while Blake's captions might be fragmentary or use informal shorthand. Feature representation differs because Casey might express the same concept in different languages, using entirely different words and grammar structures to communicate the same idea.

This natural heterogeneity reflects how people actually use technology in their daily lives. People have different communication styles, preferences, schedules, and needs. Some people use their devices constantly while others only occasionally. Some focus on text while others prefer images, videos, or audio. Some communicate in a single language while others use multiple languages. Some follow predictable patterns while others are more random in their usage. These differences in real-world data create significant challenges for federated learning systems. The system must somehow build a model that works well for everyone despite never seeing a complete picture of how different people use their devices. The training process becomes more complex because updates from different devices might conflict with each other based on their unique data patterns. Despite these challenges, addressing data heterogeneity is essential for creating effective machine learning systems that respect privacy while still providing useful features to diverse users with diverse needs.As demonstrated these variations in data distribution can cause significant performance degradation in federated visual classification tasks, with accuracy drops of up to 55% in highly heterogeneous scenarios [3].

### 2.2. Types and Advantages of Data Heterogeneity

While data heterogeneity creates many challenges for federated learning systems, it's important to understand that this diversity isn't entirely negative. In fact, when properly managed, this natural variation can actually lead to stronger, more versatile learning systems. By understanding the different types of heterogeneity that exist in real-world data,

researchers and developers can create better strategies to handle these variations and even leverage them as advantages. Rather than seeing heterogeneity as a handicap, modern approaches increasingly view it as a reflection of the real-world complexity that good machine learning systems need to handle.

## 2.3. Types of Data Heterogeneity

Data heterogeneity doesn't manifest in just one way, researchers identified that the accuracy reduction in federated learning is strongly correlated with the degree of data non-IIDness, showing that earth mover's distance between local and global distributions can quantify the severity of heterogeneity effects [4]. It appears in several distinct patterns across federated learning systems. Each type of heterogeneity presents unique challenges but also opportunities for creating more robust models. Understanding these different forms helps in designing better systems that can adapt to the natural variation in how people use technology.

## 2.4. Feature Distribution Skew

Feature distribution skew occurs when the raw data characteristics vary significantly between different clients or devices. Think of smartphone users in Japan versus those in Brazil - they likely use different languages, have different daily routines, and interact with different types of content. A Japanese user might type primarily in Japanese characters with occasional English, while a Brazilian user might use Portuguese with different slang expressions and cultural references. Even the time patterns of usage might differ due to different work cultures and time zones.

This variation extends beyond just language. Weather app users in tropical regions experience different weather patterns than those in arctic regions, creating entirely different data distributions. Financial app users in different economic systems might track completely different types of transactions or currencies. Health app users of different ages track different metrics and have different normal ranges. All these variations mean that the input features - the raw data points that the model uses for learning - can look dramatically different from one device to another, even though they're all using the same application.Recent research demonstrates that feature distribution skew can be effectively mitigated through advanced representation alignment techniques that enforce consistency between local and global feature spaces, significantly reducing performance gaps on benchmark datasets [5].

## 3. Label Distribution Skew

Label distribution skew refers to differences in how often certain outcomes or categories appear across different clients. Consider a health application that tracks patient symptoms across different hospitals. A children's hospital might see more cases of childhood diseases like chickenpox, while a cancer treatment center would record more oncology-related symptoms. A rural hospital might see more farming injuries, while an urban hospital might treat more cases related to public transportation accidents.

This imbalance means that when these institutions participate in federated learning, their local models might prioritize very different conditions. The children's hospital's model would become very good at identifying childhood illnesses but might rarely encounter heart disease symptoms. The cancer center's model would excel at cancer-related symptoms but might rarely see infectious diseases. When these models share updates, the differing frequencies of various medical conditions across institutions create a skew in which conditions the global model learns to recognize accurately. This same pattern appears in many domains - the categories or labels that are common for some users might be rare for others, creating an uneven learning environment.Strategic data augmentation and re-weighting techniques can help address label distribution skew by virtually balancing the frequency of underrepresented classes across clients, creating more equitable learning conditions without requiring data sharing. [6].

## 4. Concept Shift

Concept shift represents perhaps the most subtle but challenging form of heterogeneity. This occurs when the same input has different meanings or associations in different contexts. For example, the word "hot" in a weather app means high temperature, but in a photo-sharing app, it might indicate popularity or attractiveness. The phrase "I'm running" might indicate exercise to a fitness app but could mean lateness in a messaging app.

Cultural and regional differences amplify concept shift. A thumbs-up gesture is positive in many Western countries but can be offensive in some Middle Eastern contexts. The color red signifies danger in some cultures but celebration and good fortune in others. Even emoji usage develops different meanings across different communities and age groups. These shifting relationships between features (inputs) and labels (meanings) complicate model training because the

same pattern might need to be interpreted differently depending on context. A global model must somehow reconcile these different interpretations without access to the explicit cultural context that humans would use to disambiguate meanings.

## 5. Quantity Skew

Quantity skew happens when some clients have significantly more data than others, creating an imbalance in how much different devices contribute to the learning process. Some smartphone users might use certain apps dozens of times daily, generating thousands of data points, while others might open the same app only occasionally, creating just a few data points per month.

This imbalance appears everywhere in real-world systems. Active social media users might generate gigabytes of interaction data, while occasional users produce minimal data. People with chronic health conditions might record symptoms daily in health apps, while healthy users rarely enter information. Urban areas typically have more users than rural areas, creating geographical quantity skews. These differences in data volume mean that some devices have much richer, more comprehensive datasets for local training than others. Without proper balancing techniques, the global model might end up primarily learning from and optimizing for the most active users, potentially neglecting the needs and patterns of less frequent users who nonetheless represent important use cases.

### 5.1. Advantages of Working with Heterogeneous Data

While data heterogeneity creates numerous challenges for federated learning systems, it's not simply an obstacle to overcome. When properly understood and managed, heterogeneity can actually become a source of strength for machine learning models. The natural variations in real-world data reflect the diversity of actual users and use cases that any successful system must eventually serve. Researcher demonstrated that by properly addressing the objective inconsistency problem inherent in heterogeneous settings, federated optimization can actually leverage client diversity to achieve better overall performance than homogeneous approaches [7]. By embracing this diversity rather than fighting against it, developers can create more robust, versatile, and effective learning systems.

Data diversity leads to improved model generalization, which means the system can perform well in new, previously unseen situations. Think of this like learning a language by talking to people from many different regions rather than just one small town. If you only learn English from people in one neighborhood in Boston, you might struggle to understand English speakers from Texas, California, or England. Similarly, a model trained on homogeneous data might perform perfectly within that narrow context but fail when encountering slightly different scenarios. Models that learn from diverse data sources gain exposure to a wider range of patterns, edge cases, and variations. This broader exposure helps them develop more flexible and adaptable internal representations that can handle new situations more gracefully. For example, a speech recognition system trained on diverse accents, background noise conditions, and speaking styles will be much more likely to correctly understand a new user with an unfamiliar accent than a system trained on a single accent in quiet laboratory conditions.

Heterogeneous data creates valuable opportunities for personalization that wouldn't exist with uniform data. The variations between users become information that the system can leverage to tailor experiences. Rather than forcing everyone into the same one-size-fits-all model, federated learning can identify patterns in how different groups of users behave and adapt accordingly. A text prediction system might recognize that some users prefer formal language while others use more casual expressions, allowing it to offer different suggestions to different people. A fitness app could recognize that exercise patterns vary dramatically between users of different ages, fitness levels, or geographic regions, and adjust its recommendations accordingly. Without heterogeneity, these personalization opportunities would be much more limited. The differences between users, which initially appear as a challenge, ultimately become the foundation for creating more relevant and useful experiences tailored to individual needs and preferences.

Systems designed to handle heterogeneous data tend to have greater real-world applicability because they're built from the beginning to cope with the messiness and complexity of actual usage. The real world doesn't provide clean, balanced, perfectly labeled datasets - it provides chaotic, imbalanced, noisy information that varies unpredictably. By developing techniques that work under these challenging conditions, federated learning researchers create systems that don't just work in laboratory settings but can function effectively in actual deployments. A recommendation system that can handle users with vastly different preferences and usage patterns will transition more smoothly from development to production than one optimized for an artificially uniform dataset. The robustness required to handle heterogeneity becomes a competitive advantage when systems move beyond controlled environments into the unpredictable diversity of real-world applications.

Working with heterogeneous data often requires techniques that inherently enhance privacy protection. Because the system must be designed from the ground up to handle variations between users without direct access to their raw data, it tends to develop approaches that are less dependent on specific individual information. Differential privacy techniques, which add careful noise to protect individual data points, work particularly well in heterogeneous environments where the system already expects variation. Federated learning systems built for heterogeneous data often use aggregation methods that further obscure individual contributions, such as secure multi-party computation or homomorphic encryption. These privacy-enhancing techniques become natural extensions of the approaches needed to handle heterogeneity effectively, creating systems that are both more versatile and more respectful of user privacy.

Finally, systems designed for heterogeneous environments tend to develop better fault tolerance because they're already built to handle unusual or unexpected patterns. When a system expects uniformity, any deviation might cause catastrophic failures. But when a system is designed around the premise that inputs will vary widely, it develops internal mechanisms to gracefully handle outliers, anomalies, and edge cases. This resilience means that federated learning systems optimized for heterogeneous data can often continue functioning effectively even when encountering corrupted data, unusual usage patterns, or partial system failures. For example, a content recommendation system trained on diverse user behaviors would be less likely to completely break down when encountering a user with unusual preferences - it would simply recognize this as another form of the heterogeneity it was designed to accommodate. This natural robustness represents a significant advantage in real-world deployments where unpredictable events and unusual patterns inevitably occur.

## 5.2. Challenges of Data Heterogeneity in Federated Learning

While data heterogeneity offers some advantages, it also introduces significant hurdles that make federated learning more complex than traditional centralized approaches. These challenges can undermine the effectiveness, efficiency, and fairness of federated systems if not properly addressed. Understanding these obstacles is crucial for developing strategies to mitigate their effects and build federated learning systems that work well despite the natural variations in real-world data. As comprehensively reviewed, data heterogeneity remains one of the most fundamental challenges in federated learning, affecting every aspect from optimization convergence to fairness and system efficiency [8].

## 5.3. Convergence Issues

One of the most fundamental challenges of data heterogeneity involves the learning process itself. When a federated system attempts to combine model updates from diverse data sources, the training process often becomes unstable and difficult to manage. This is similar to trying to navigate a ship when receiving conflicting directions from multiple captains - the vessel might move erratically or struggle to make progress toward its destination.

Slow convergence often plagues heterogeneous federated learning systems. The training process typically requires many more rounds to reach a good solution compared to homogeneous data settings. Imagine teaching a class where each student has vastly different background knowledge - you would need to spend more time ensuring everyone understands the material. Similarly, when client devices have different data distributions, the global model needs more updates to adequately learn from all the diverse patterns. A model might quickly become accurate for the most common data patterns but require many additional rounds to perform well across all the variations.

Oscillation represents another common convergence problem. The model parameters might bounce back and forth between different values without settling on a good solution. This happens because updates from different client devices might push the model in contradictory directions. For example, a text prediction model might receive updates from business users suggesting formal language predictions, only to then receive updates from casual users suggesting informal language. The global model might swing between these different styles without finding a balanced solution that works well for everyone. These oscillations waste computational resources and can prevent the model from ever reaching optimal performance.

Perhaps most concerning is premature convergence to suboptimal solutions. The system might appear to stabilize, giving the impression that training is complete, but actually settle on a solution that works moderately well for most clients without excelling for any of them. This is like finding a compromise that nobody is particularly happy with. The global model might reach a local minimum that balances the competing needs of different data distributions without discovering a more creative solution that could better serve all users. This premature convergence can be difficult to detect since the model appears stable, but its performance remains disappointing compared to what might be possible with better optimization techniques. Theoretical analysis proves that the convergence rate of FedAvg degrades proportionally to the degree of data heterogeneity, requiring more careful selection of learning rates and communication rounds in non-IID settings [9].

## 5.4. Model Bias and Fairness Concerns

Heterogeneous data introduces serious risks of bias and fairness problems in federated learning systems. When client data varies significantly, the global model often ends up favoring the majority patterns while performing poorly for minority groups or edge cases. This bias occurs naturally unless specific countermeasures are implemented.

The problem stems from how federated averaging typically works. Clients with more common data patterns tend to have more influence on the final model, especially if there are more such clients or if they have larger datasets. For example, if a speech recognition system is trained across users from many countries, but 70% of participants are from English-speaking regions, the resulting model will likely perform much better for English speakers than for speakers of less represented languages. This creates an inherently unfair system where certain users receive higher quality service than others based solely on whether their usage patterns match the majority.

These biases can reinforce existing social inequalities when they align with demographic differences. If certain demographic groups use technology in distinctive ways or have unique needs, a biased federated learning system might systematically underserve these populations. A health monitoring app might become excellent at detecting symptoms common in the majority population but miss critical indicators that present differently in minority groups. A language model might handle standard dialects well but struggle with regional variations or accents. These disparities raise serious ethical concerns about whether federated learning systems are providing equitable service to all users.

The problem becomes especially challenging because the federated nature of the system makes bias harder to detect and address. In a centralized system, researchers could directly analyze the training data to identify underrepresented groups and implement targeted solutions. But in federated learning, the raw data remains invisible to the central server, making it difficult to even measure bias, let alone correct it. This opacity creates additional barriers to ensuring fairness in heterogeneous federated systems.Fairness-aware federated learning approaches that explicitly incorporate equity objectives into the learning process can help counteract these systemic biases, ensuring more balanced performance across diverse subpopulations while still maintaining privacy guarantees [8].

## 5.5. Communication Inefficiency

Heterogeneous data distributions significantly increase the communication burden in federated learning systems. When client data varies widely, the system typically requires many more communication rounds to reach a good solution. Each round involves sending model updates from client devices to the central server and distributing new global model versions back to the clients. This increased communication has real costs in terms of bandwidth usage, energy consumption, and time.

For mobile devices, these costs are particularly concerning. Frequent model updates consume battery power and data allowances. A smartphone participating in a federated learning system might drain its battery faster or use up the user's monthly data plan. In regions with limited connectivity or expensive data plans, this increased communication burden could effectively exclude many potential participants from the system. Even in areas with good connectivity, the additional network traffic from millions of devices sending frequent model updates could place significant strain on infrastructure.

The communication problem compounds when we consider that heterogeneous data might require larger or more complex models to capture all the variations. A simple model might adequately handle homogeneous data, but representing the diverse patterns in heterogeneous data often demands more parameters. These larger models require more bandwidth to transfer, further increasing the communication costs. The combination of more frequent updates and larger model sizes can make heterogeneous federated learning prohibitively expensive in terms of communication resources.

## 5.6. Privacy-Utility Trade-offs

While federated learning inherently enhances privacy by keeping raw data on local devices, heterogeneous data introduces additional privacy challenges that often result in difficult tradeoffs between privacy protection and model utility. The unique or unusual patterns in heterogeneous data can sometimes be more identifiable, potentially increasing privacy risks.

When a system encounters highly distinctive data patterns from a small subset of users, it faces a dilemma. Including these patterns helps the model serve those users better, but it might also make their contributions more recognizable in the final model. For example, if only a small number of users write in a particular language or dialect, their model

updates might be distinctive enough that an adversary could potentially identify information about them from the global model. This creates a situation where improving service for unique user groups might simultaneously increase their privacy risks.

To address these concerns, federated learning systems often implement additional privacy protections such as differential privacy techniques, which add carefully calibrated noise to mask individual contributions. However, these privacy mechanisms typically reduce model accuracy, especially for minority patterns that are already challenging to learn. The more noise added to protect privacy, the harder it becomes for the model to pick up on subtle or uncommon patterns in the data. This creates a direct tradeoff - stronger privacy protections often mean worse performance for users with non-mainstream data patterns, the very users who might already be underserved by the model.

## 5.7. Computational Imbalance

Data heterogeneity often creates significant imbalances in the computational demands placed on different client devices. Some clients might need to process much larger or more complex datasets than others, creating an unfair distribution of the computational burden.

This imbalance can manifest in several ways. Clients with more data naturally require more processing time and energy to complete their local training. A smartphone with years of user data might need significantly more computation than a new device with minimal history. Similarly, certain types of data require more intensive processing - video or audio data typically demands more computation than simple text or numerical data. If some users primarily generate computationally intensive data types while others generate simpler data, the computational demands will be unevenly distributed.

The consequences of this imbalance can be serious for resource-constrained devices. Older smartphones, IoT devices with limited processing power, or devices with battery limitations might be unable to participate effectively in federated learning if the computational demands are too high. This could create a system that systematically excludes certain devices or users, potentially introducing further biases into the learning process. For example, if only high-end devices can practically participate, the resulting model might be optimized for users who can afford such devices, neglecting the needs and patterns of users with more basic technology.

This computational imbalance also creates practical challenges for system design. Should the federated learning system wait for all clients to complete their local training, potentially creating long delays as it waits for the slowest devices? Or should it proceed with only the faster clients, potentially introducing selection bias? These questions highlight how computational imbalance creates difficult tradeoffs between inclusivity, efficiency, and fairness in heterogeneous federated learning systems.

## 5.8. Current Solutions and Future Plans

The challenges of data heterogeneity in federated learning have sparked creative solutions from researchers around the world. Rather than viewing heterogeneity as merely an obstacle, the field has begun to develop innovative approaches that not only mitigate its negative effects but sometimes even leverage diversity as a strength. These solutions range from fundamental changes in how models are structured to sophisticated techniques for combining and managing updates from diverse sources. Wang et al. demonstrated that reinforcement learning techniques can dynamically optimize federated learning processes on non-IID data, automatically selecting the best clients for each round based on their expected contributions to model improvement [10]. As the field continues to mature, emerging research directions promise even more effective ways to handle the natural variations in real-world data.

## 5.9. Current Solutions

Personalized federated learning represents one of the most promising approaches to handling heterogeneous data. This strategy acknowledges that a single model cannot optimally serve all users and instead builds systems that adapt to individual data patterns. Rather than forcing every client to conform to the same global model, personalized approaches create flexible models that can adjust to each user's unique needs. Meta-learning techniques train models that are specifically designed to adapt quickly to new patterns with minimal data, making them well-suited for personalization. For example, a smartphone keyboard might start with a global language model but rapidly customize its predictions based on a user's unique writing style. Local fine-tuning allows each device to make final adjustments to the global model based on its specific data, creating a semi-personalized experience. Some advanced systems even maintain both global and personal components, using the global model for common patterns while personal layers handle user-specific behaviors. This personalization approach recognizes heterogeneity as a natural feature rather than a bug and designs

systems that embrace these differences.Fallah et al. established theoretical guarantees for personalized federated learning using model-agnostic meta-learning, proving that it can achieve near-optimal performance even with highly heterogeneous data by learning an initialization that allows quick adaptation to individual clients [11].

Robust aggregation methods help create more stable and reliable global models despite the varying updates coming from heterogeneous clients. Traditional averaging techniques can be easily skewed by outliers or conflicting updates, but advanced aggregation methods filter and combine client contributions more intelligently. Median-based aggregation replaces simple averaging with median values, reducing the impact of extreme outliers that might come from unusual data distributions. Trimmed mean approaches discard some percentage of the most extreme updates before averaging, creating a more balanced and representative combination. Other techniques weight contributions based on the quality or reliability of updates, giving more influence to clients that provide consistent, generalizable improvements. These robust aggregation methods create more stable convergence patterns and help the global model find better solutions that work across diverse data distributions. By intelligently combining perspectives from different clients, these approaches turn the potential chaos of heterogeneous updates into a strength that improves model quality. Researchers proposed an agnostic federated learning approach that optimizes for the worst-case client distribution rather than the average, resulting in more robust models that perform equitably across heterogeneous clients without requiring explicit fairness constraints [12].

Federated distillation offers a way to create smaller, more efficient personalized models by transferring knowledge from larger, more complex models. Rather than sharing the full model structure and parameters, distillation compresses the essential patterns and behaviors into a more compact form. A large global teacher model learns from the collective knowledge of all clients, capturing patterns across the full diversity of data. Then, smaller student models are created for individual devices, containing only the knowledge most relevant to each user's needs. This approach drastically reduces the communication costs since smaller models require less bandwidth to transfer. It also decreases the computational demands on client devices, making participation more accessible to resource-constrained hardware. Federated distillation provides a particularly elegant solution to handling heterogeneity because it naturally separates universal knowledge that applies to everyone from specialized knowledge that only matters to specific users or use cases. The result is a system that efficiently distributes intelligence across the network while respecting the unique needs of individual clients.Unlike traditional federated learning that exchanges model parameters whose size depends on the model architecture, federated distillation exchanges only output logits that depend on the number of classes, dramatically reducing communication overhead especially for deep neural networks. This communication efficiency enables federated distillation to achieve comparable accuracy to federated learning while requiring up to 26 times less data transmission, making it particularly suitable for resource-constrained mobile environments where bandwidth and energy consumption are critical constraints [13].

Client clustering takes a divide-and-conquer approach to heterogeneity by grouping similar clients together before training. Instead of forcing wildly different data distributions into a single model, clustering techniques identify natural groupings of clients based on their data characteristics. Each cluster then trains its own federated model, creating several specialized models rather than one universal model. For example, a language learning app might create separate clusters for beginners, intermediate, and advanced learners, with each group training models specifically optimized for their skill level. Clustering can be based on various factors, including data distributions, usage patterns, geographic location, or device types. This approach significantly reduces the negative impacts of extreme heterogeneity by ensuring that each model only needs to handle a more manageable level of variation. Client clustering effectively transforms a highly heterogeneous global problem into several more homogeneous local problems that are easier to solve. The tradeoff is increased system complexity, as the federated learning system must now maintain multiple models and determine which one to apply for each client. However, clustering methods must adapt to data drifts that naturally occur when clients' data distribution changes over time, as static clustering can soon become as heterogeneous as having no clustering at all. Fielding addresses this challenge by combining per-client migrations with selective global re-clustering, balancing the need to maintain optimal clusters against the instability introduced by frequent re-clustering [14].

Adaptive client participation strategies intelligently select which clients should participate in each training round based on their data characteristics and potential contributions. Rather than randomly choosing devices or including everyone in every round, these approaches make strategic decisions about participation to improve efficiency and fairness. Some systems prioritize clients with data that is currently underrepresented in the model, ensuring that minority patterns receive adequate attention. Others might temporarily exclude clients with extremely unusual data to stabilize the early learning process, gradually incorporating these outliers in later rounds. Adaptive approaches might also consider practical factors like device battery levels, connection quality, or computational capacity when deciding participation. By thoughtfully orchestrating which clients contribute when, these methods improve the learning dynamics in

heterogeneous environments. They can reduce the number of communication rounds needed while simultaneously ensuring that the final model better represents the full diversity of client data.

## 6. Future Research Directions

The field of federated learning continues to evolve with several promising research directions that address the challenges of data heterogeneity.

### 6.1. Advanced Heterogeneity Metrics

The development of advanced heterogeneity metrics represents a crucial frontier in federated learning research. Currently, researchers lack standardized, comprehensive ways to measure and characterize exactly how data differs across clients. Future work aims to create sophisticated tools that can quantify different aspects of heterogeneity without requiring access to the raw data itself. These metrics might measure statistical distances between client distributions, identify specific types of variation, or detect problematic patterns that could destabilize training. Better measurement tools would enable more targeted solutions by helping systems understand precisely what kind of heterogeneity they're dealing with. For example, knowing whether variation comes primarily from feature distribution differences versus label distribution differences would help select the most appropriate mitigation strategies. Advanced metrics would also help in monitoring and evaluating system performance across different subgroups, ensuring that improvements benefit all users equitably. As federated learning scales to even more diverse environments, these measurement tools will become increasingly important for understanding system behavior and identifying opportunities for improvement.

### 6.2. Cross-Silo Federated Learning

Cross-silo federated learning focuses on collaboration between organizations rather than individual devices, presenting unique heterogeneity challenges that require specialized solutions. When hospitals, banks, or government agencies participate in federated learning, their data distributions often differ dramatically due to their distinct client bases, services, and operational patterns. These differences tend to be more systematic and extreme than the variations between individual users' devices. Future research aims to develop techniques specifically designed for this organizational context, where each participant may have substantial computational resources but highly specialized data. This might include specialized privacy-preserving techniques for sensitive institutional data, contractual frameworks for fair collaboration between competing organizations, and methods to handle the structured heterogeneity that emerges from different organizational practices. As more industries recognize the value of collaborative learning without data sharing, cross-silo federated learning will become increasingly important for applications like medical research, financial risk modeling, and public sector services, all contexts where data heterogeneity presents particularly complex challenges.

### 6.3. Dynamic Adaptation Mechanisms

Dynamic adaptation mechanisms represent another promising research direction, focusing on systems that automatically adjust to changing data distributions over time. Most current federated learning approaches assume relatively stable client characteristics, but real-world data continuously evolves as user behaviors change, new devices enter the system, and external circumstances shift. Future systems will need to detect these changes automatically and adapt their strategies accordingly. This might involve continuously monitoring distribution shifts, automatically adjusting personalization parameters, or dynamically reconfiguring client clusters as patterns evolve. For example, a model might detect seasonal changes in user behavior and preemptively adjust its learning approach, or it might recognize when a previously unusual data pattern is becoming more common and increase its representation in the global model. These adaptive systems would be particularly valuable in domains with rapidly changing conditions, such as health monitoring during disease outbreaks, consumer behavior during economic shifts, or transportation patterns during infrastructure changes.

### 6.4. Hardware-Aware Federated Learning

Hardware-aware federated learning recognizes that data heterogeneity often correlates with device heterogeneity, requiring integrated solutions that address both challenges simultaneously. Future research will increasingly consider how varying device capabilities — processing power, memory, battery capacity, connectivity — interact with data variations to influence system performance. This integrated approach might involve automatically tailoring local training processes to device capabilities, prioritizing critical model components for resource-constrained devices, or dynamically adjusting participation based on both data relevance and device status. For example, systems might send smaller, specialized model segments to devices with limited memory while providing more comprehensive models to

high-capacity devices. This research direction acknowledges that in real-world deployments, data and device characteristics are deeply intertwined, and effective solutions must address this complexity holistically. As federated learning expands to include an even wider range of devices, from powerful servers to tiny IoT sensors, hardware-aware approaches will become increasingly essential for creating systems that work well across the full spectrum of participants.

## 6.5. Strengthening Theoretical Foundations

Stronger theoretical foundations will provide crucial insights into how heterogeneous data fundamentally affects model convergence, generalization, and performance. While practical solutions have advanced rapidly, the mathematical understanding of federated learning with non-IID data remains incomplete. Future research will develop more comprehensive theoretical models that can predict system behavior, establish performance guarantees, and guide algorithm design in heterogeneous environments. This theoretical work might explore concepts like convergence bounds under different types of heterogeneity, optimal aggregation strategies for specific distribution patterns, or fundamental limits on what can be achieved with various privacy constraints. Better theory will enable more principled system design, moving beyond heuristic approaches to solutions with provable properties and guarantees. This research direction represents the deep scientific work needed to transform federated learning from a collection of effective techniques into a rigorously understood field with predictable behaviors and outcomes.

## 6.6. Developing Standardized Benchmarks

Standardized benchmarks will help researchers compare different approaches and measure progress in addressing heterogeneity challenges. Currently, federated learning research often uses different datasets, heterogeneity models, and evaluation metrics, making it difficult to directly compare competing methods. Future work will establish common test scenarios that realistically represent the types of heterogeneity encountered in different domains. These benchmarks might include synthetic datasets with controllable heterogeneity parameters, carefully curated real-world datasets that capture typical variation patterns, or simulation environments that model both data and system dynamics. Domain-specific benchmarks will be particularly valuable, recognizing that heterogeneity manifests differently in healthcare, finance, mobile applications, and other contexts. By evaluating different methods on the same challenging scenarios, these benchmarks will accelerate progress by clearly identifying which approaches work best for specific heterogeneity challenges. They will also help bridge the gap between academic research and practical applications by ensuring that solutions address realistic problems rather than simplified abstractions.

## 7. Conclusion

Data heterogeneity represents both a significant challenge and an opportunity for federated learning. While the non-uniform nature of real-world data complicates the training process, addressing these challenges drives innovation in machine learning algorithms, personalization techniques, and privacy-preserving methods. As federated learning matures, its ability to effectively handle heterogeneous data will determine its success in critical applications like healthcare, finance, and smart devices. The future of federated learning lies not in eliminating heterogeneity, but in embracing it as a natural characteristic of decentralized systems and developing methods that thrive in diverse data environments.

## References

[1] H. Brendan McMahan, et al., "Communication-Efficient Learning of Deep Networks from Decentralized Data," Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS), 2023. [Online]. Available: https://arxiv.org/pdf/1602.05629

[2] Hangyu Zhu, et al., "Federated Learning on Non-IID Data: A Survey," Elsevier, 2021. [Online]. Available: https://arxiv.org/pdf/2106.06843

[3] Tzu-Ming Harry Hsu, et al., "Measuring the Effects of Non-Identical Data Distribution for Federated Visual Classification," arXiv preprint arXiv:1909.06335, 2019. [Online]. Available: https://arxiv.org/pdf/1909.06335

[4] Yue Zhao, et al., "Federated Learning with Non-IID Data," arXiv preprint arXiv:1806.00582, 2022. [Online]. Available: https://arxiv.org/pdf/1806.00582

[5] Yunlu Yan, et al., "A Simple Data Augmentation for Feature Distribution Skewed Federated Learning," arXiv:2306.09363v2 [cs.LG] 6 Dec 2024. [Online]. Available: https://arxiv.org/pdf/2306.09363

[6]     Jie Zhang, et al., "Federated Learning with Label Distribution Skew via Logits Calibration," Proceedings of the 39th International Conference on Machine Learning, PMLR 162:26311-26329, 2022.. [Online]. Available: https://proceedings.mlr.press/v162/zhang22p.html

[7]     Jianyu Wang, et al., "Tackling the Objective Inconsistency Problem in Heterogeneous Federated Optimization," 34th Conference on Neural Information Processing Systems (NeurIPS 2020),. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2020/file/564127c03caab942e503ee6f810f54fd-Paper.pdf

[8]     Peter Kairouz, et al., "Advances and Open Problems in Federated Learning," Foundations and Trends in Machine Learning, vol. 14, no. 1–2, pp. 1-210, 2021. [Online]. Available: https://arxiv.org/pdf/1912.04977

[9]     Xiang Li, et al., "On The Convergence Of Fedavg On Non-Iid Data," in International Conference on Learning Representations (ICLR), 2020. [Online]. Available: https://arxiv.org/pdf/1907.02189

[10]    Hao Wang, et al., "Optimizing Federated Learning on Non-IID Data with Reinforcement Learning," in IEEE INFOCOM 2020 - IEEE Conference on Computer Communications, 2020, pp. 1698-1707. [Online]. Available: https://iqua.ece.toronto.edu/papers/hwang-infocom20.pdf

[11]    Alireza Fallah, et al., "Personalized Federated Learning with Theoretical Guarantees: A Model-Agnostic Meta-Learning Approach," 34th Conference on Neural Information Processing Systems (NeurIPS 2020). [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2020/file/24389bfe4fe2eba8bf9aa9203a44cdad-Paper.pdf

[12]    Mehryar Mohri, et al., "Agnostic Federated Learning," in Proceedings of the 36th International Conference on Machine Learning, 2019, pp. 4615-4625. [Online]. Available: https://arxiv.org/pdf/1902.00146

[13]    Eunjeong Jeong, et al., "Communication-Efficient On-Device Machine Learning: Federated Distillation and Augmentation under Non-IID Private Data", arXiv:1811.11479v2 [cs.LG] 19 Oct 2023. [Online]. https://arxiv.org/pdf/1811.11479

[14]    Minghao Li, et al., "Federated Learning Clients Clustering with Adaptation to Data Drifts", arXiv:2411.01580v1 [cs.LG] 3 Nov 2024. [Online] https://arxiv.org/pdf/2411.01580