

ReSoFed: Reliability-Guided Model Souping for Robust Federated Learning in Heterogeneous Classroom Environments

Anonymous CVPR submission

Paper ID ****

Abstract

001 Computer vision systems are increasingly used in educa-
002 tional environments to analyze classroom activities and
003 support data-driven learning analytics. However, class-
004 room visual data often contain sensitive student informa-
005 tion and cannot be centralized across institutions due to pri-
006 vacy and governance constraints. Federated learning (FL)
007 enables collaborative model training without sharing raw
008 data, but its performance can degrade when visual data
009 originate from heterogeneous capture environments. Un-
010 der such conditions, unreliable client updates can introduce
011 negative transfer during model aggregation. In this work,
012 we propose **ReSoFed**¹, a privacy-preserving, reliability-
013 guided federated aggregation framework designed to im-
014 prove robustness under heterogeneous classroom environ-
015 ments. The framework estimates the cross-domain gener-
016 alization capability of client models using a heterogeneous
017 server-held validation set and incorporates this signal into
018 a two-stage aggregation process. Specifically, a greedy
019 model-soup procedure identifies a subset of reliable client
020 models whose weight-space combination improves valida-
021 tion performance, followed by reliability-aware weighted
022 aggregation. Experiments on the SCB dataset across multi-
023 ple CNN and transformer backbones demonstrate that Re-
024 SoFed consistently outperforms standard federated learn-
025 ing baselines under heterogeneous visual conditions while
026 preserving data privacy.

027 1. Introduction

028 Computer vision technologies are increasingly integrated
029 into educational environments to support classroom analyt-
030 ics, student engagement monitoring, and automated analy-
031 sis of learning behaviors. By analyzing visual cues such as
032 classroom activities, participation patterns, and student in-
033 teractions, these systems can provide valuable insights for
034 educators [24]. Recent research has also explored multi-

modal approaches that combine visual, behavioral, and con- 035
textual signals to detect collaborative learning states and 036
support AI-assisted classroom analytics [2, 4]. However, 037
deploying such vision-based educational systems across in- 038
stitutions presents a fundamental challenge: educational 039
data is highly sensitive and subject to strict privacy regu- 040
lations and institutional policies. Video streams, classroom 041
recordings, and student behavior data cannot typically be 042
centralized or shared across institutions due to ethical and 043
privacy constraints. 044

Federated learning (FL) offers a promising solution by 045
enabling collaborative model training across distributed 046
institutions without exchanging raw data [11]. In this 047
paradigm, each institution trains a model locally on its pri- 048
vate data and shares only model parameters with a central 049
server for aggregation [14], allowing knowledge sharing 050
while preserving privacy. However, applying FL to real- 051
world classroom vision systems introduces additional chal- 052
lenges. Classroom environments differ significantly across 053
institutions in terms of camera hardware, lighting condi- 054
tions, viewpoints, and recording setups, leading to substan- 055
tial distribution shifts in visual data [7]. These variations 056
produce highly non-IID client distributions [1], which can 057
degrade the performance and stability of federated learning 058
models and limit their ability to generalize across heteroge- 059
neous institutions. 060

A critical but often overlooked factor contributing to 061
this issue is that not all client updates are equally reli- 062
able. Clients operating under degraded visual conditions 063
may produce models that overfit to environment-specific ar- 064
tifacts. Conventional federated aggregation strategies, such 065
as Federated Averaging (FedAvg) [18], implicitly assume 066
that all client updates contribute equally useful informa- 067
tion and therefore aggregate models using uniform or data- 068
size-weighted averaging. Under heterogeneous visual con- 069
ditions, this assumption can lead to negative transfer, where 070
poorly generalizing client updates degrade the performance 071
of the global model [22]. This challenge becomes particu- 072
larly pronounced in educational settings, where classroom 073
capture conditions can vary substantially across institutions. 074

¹Codes and data are available at [TBA](#)

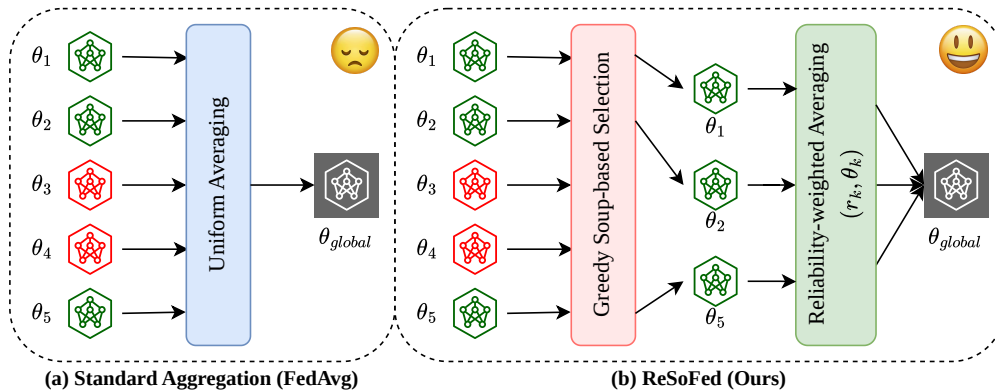


Figure 1. Comparison between (a) standard federated aggregation (FedAvg) and (b) the proposed **ReSoFed** framework. FedAvg uniformly aggregates all client updates ($\theta_1, \dots, \theta_5$), which may introduce negative transfer under heterogeneous visual conditions. In contrast, ReSoFed performs greedy model-soup-based client selection to filter unreliable clients and then applies reliability-aware weighted aggregation to produce the global model θ_{global} . Strongly generalizing clients are shown in **green**, while weak clients are marked in **red**.

075 In this work, we argue that cross-domain reliability can
 076 serve as a key signal for federated aggregation in privacy-
 077 preserving educational vision systems. Client models that
 078 generalize well beyond their local classroom conditions are
 079 more reliable for constructing a robust global model. Mo-
 080 tivated by this insight, we propose **ReSoFed (Reliability-**
 081 **guided Soup for Federated Learning)**, a federated aggrega-
 082 tion framework that explicitly estimates and leverages client
 083 reliability during model aggregation. Instead of assuming
 084 uniform contributions from all clients, ReSoFed evaluates
 085 locally trained models on a heterogeneous server-held vali-
 086 dation set to estimate their cross-domain generalization ca-
 087 pability. Building on this signal, the framework performs
 088 a two-stage aggregation process. First, a greedy model-
 089 soup-based client selection procedure iteratively combines
 090 candidate models only when their inclusion improves vali-
 091 dation performance, filtering out poorly generalizing up-
 092 dates. Second, reliability-aware weighted aggregation is ap-
 093 plied over the selected clients, assigning greater influence
 094 to models with stronger cross-domain generalization. Fig-
 095 ure 1 provides a high-level comparison between the stan-
 096 dard FedAvg aggregation and the proposed ReSoFed frame-
 097 work. We evaluate ReSoFed on the SCB dataset [30] for
 098 classroom activity recognition, where client heterogeneity
 099 is simulated through diverse augmentations. Experimental
 100 results show that ReSoFed mitigates negative transfer
 101 while preserving privacy, requires no modification to local
 102 training objectives, and remains compatible with diverse ar-
 103 chitectures, including both convolutional and transformer-
 104 based models.

105 **Our major contributions are as follows:**

- 106 • We introduce **ReSoFed**, a privacy-preserving, reliability-
 107 guided federated aggregation framework that models
 108 client reliability via cross-domain validation performance
 109 under heterogeneous visual conditions.

- We propose a cross-domain reliability estimation strategy that evaluates client models on a heterogeneous server-held validation set to assess their generalization capability beyond local data distributions.
- A two-stage reliability-guided aggregation mechanism is developed that integrates greedy model-soup-based client selection with reliability-aware weighted averaging, enabling the server to filter weakly generalizing updates and prioritize robust client models during global aggregation.
- Extensive experiments demonstrate that ReSoFed consistently improves robustness and global model performance over standard federated learning baselines while remaining privacy-preserving, architecture-agnostic, and computationally lightweight.

2. Related Works

2.1. Federated Learning

Federated Learning (FL) is a distributed machine learning paradigm that enables multiple clients to collaboratively train a shared model without exchanging raw data, thereby preserving data privacy and reducing the need for centralized data collection [11, 18]. In a typical FL setup, each client performs local optimization on its private dataset and sends model weights to a central server, which aggregates them to form a global model. The most widely used baseline is Federated Averaging (FedAvg) [18], which aggregates client updates through data-size-weighted averaging. Despite its advantages, FL faces several challenges that are not present in traditional distributed learning settings. In particular, data heterogeneity across clients, often referred to as non-IID distribution, can significantly degrade model convergence and generalization performance [7, 14]. Such heterogeneity commonly arises in real-world deployments due to variations in device characteristics, environmental

143 conditions, or user behavior. To address the challenges
144 posed by heterogeneous data, several extensions to FedAvg
145 have been proposed. For example, FedProx introduces a
146 proximal regularization term in the local optimization ob-
147 jective to stabilize training under heterogeneous conditions
148 [15]. Other approaches focus on improving optimization
149 and communication efficiency through modified aggrega-
150 tion strategies or adaptive optimization schemes [10, 23].
151 These methods highlight the importance of designing ag-
152 gregation strategies that can effectively handle variability
153 across distributed clients while maintaining the privacy-
154 preserving benefits of federated learning.

155 2.2. Client Selection and Aggregation Strategies

156 Beyond the standard Federated Averaging (FedAvg) al-
157 gorithm, a large body of work has focused on improv-
158 ing the aggregation and client participation strategies in
159 federated learning. The aggregation rule plays a central
160 role in federated optimization, as it integrates knowledge
161 from distributed client models to produce a global model
162 while addressing challenges such as statistical heterogene-
163 ity and communication constraints [19, 21]. One line of
164 research focuses on improving aggregation through adap-
165 tive optimization strategies. Methods such as FedOpt in-
166 troduce server-side adaptive optimization schemes to stabilize
167 global model updates under heterogeneous client distribu-
168 tions [23]. Similarly, FedNova addresses objective incon-
169 sistency caused by varying local update steps across clients,
170 improving convergence behavior in heterogeneous settings
171 [28]. These approaches aim to improve the stability and ef-
172 ficiency of global aggregation without explicitly modeling
173 client reliability. More recently, several works have investi-
174 gated reliability-driven or contribution-aware aggregation
175 mechanisms. These methods attempt to quantify the impor-
176 tance of each client update using metrics such as local loss,
177 gradient similarity, or contribution scores, and adjust ag-
178 gregation weights accordingly [9]. While these approaches
179 demonstrate improvements under heterogeneous data con-
180 ditions, many of them require additional information from
181 clients or introduce significant computational overhead.

182 In contrast, our approach focuses on estimating cross-
183 domain reliability using a server-held validation set and in-
184 tegrating this signal into the aggregation process through
185 reliability-guided client selection and weighted model ag-
186 gregation. This design enables improved robustness under
187 heterogeneous visual conditions while maintaining compat-
188 ibility with standard federated training pipelines.

189 2.3. Computer Vision for Education

190 Computer vision has increasingly been applied to educa-
191 tional environments to analyze student behaviors, engage-
192 ment levels, and classroom interactions. Early research pri-
193 marily focused on engagement recognition using facial ex-

194 pressions, head pose, and gaze patterns captured through
195 webcams in online learning settings. Datasets such as
196 DAiSEE have enabled the development of models that es-
197 timate student engagement levels from video streams by
198 analyzing visual cues such as facial movements and atten-
199 tion patterns [5, 12]. More recent works leverage deep
200 neural networks, including convolutional and transformer-
201 based architectures, to improve engagement recognition us-
202 ing spatio-temporal visual signals [20, 25]. These studies
203 demonstrate the potential of computer vision to support in-
204 telligent tutoring systems and provide educators with in-
205 sights into student learning behavior. Another important re-
206 search direction focuses on detecting and recognizing class-
207 room activities from visual data. For example, the Student
208 Classroom Behavior (SCB) dataset provides annotated im-
209 ages of common activities such as reading, writing, rais-
210 ing hands, and listening, enabling the development of au-
211 tomated behavior recognition systems [26]. Computer vi-
212 sion models based on object detection and action recogni-
213 tion frameworks have been proposed to analyze classroom
214 scenes and identify student activities associated with en-
215 gagement and participation [3]. These approaches aim to
216 support classroom analytics and enhance teaching effective-
217 ness by providing automated behavioral monitoring.

218 Despite significant progress in federated learning, aggrega-
219 tion strategies, and computer vision for education, sev-
220 eral challenges remain when deploying vision models in
221 real-world educational environments. Federated learning
222 enables privacy-preserving collaborative training but suf-
223 fers from degraded convergence and generalization under
224 heterogeneous data distributions. Existing aggregation and
225 client selection methods improve optimization stability but
226 often do not explicitly capture cross-domain reliability of
227 client models [13]. Meanwhile, most computer vision ap-
228 proaches for educational analytics, including engagement
229 and behavior recognition, assume centralized training on
230 pooled datasets [5, 26]. In practice, educational data are
231 distributed across institutions and captured under highly
232 heterogeneous visual conditions. To address these chal-
233 lenges, we propose a reliability-guided federated learning
234 framework that combines cross-domain validation-based
235 reliability estimation with model-soup-inspired client se-
236 lection and reliability-aware aggregation, enabling robust
237 and privacy-preserving learning across heterogeneous class-
238 room environments.

239 3. Methodology

240 3.1. Problem Formulation

241 We consider a federated image classification setting with
242 K distributed clients. Each client $k \in \{1, \dots, K\}$ pos-
243 sesses a private local dataset $\mathcal{D}_k = \{(x_i^k, y_i^k)\}_{i=1}^{n_k}$ of size
244 n_k , where $x_i^k \in \mathcal{X}$ denotes an image and $y_i^k \in \{1, \dots, C\}$

245 denotes its corresponding label. In our setting, $C = 11$ for
246 classroom action recognition. Due to heterogeneous visual
247 capture conditions, the local datasets \mathcal{D}_k are non-identically
248 distributed (non-IID) across clients.

249 Let $f_\theta : \mathcal{X} \rightarrow \mathbb{R}^C$ denote a shared model parameterized
250 by $\theta \in \mathbb{R}^d$. In standard federated learning, the objective is
251 to minimize the global empirical loss:

$$252 \min_{\theta} F(\theta) = \sum_{k=1}^K \frac{n_k}{\sum_{j=1}^K n_j} F_k(\theta).$$

253 where $F_k(\theta)$ denotes the local empirical loss at client k .

254 Under heterogeneous visual conditions, direct aggrega-
255 tion of locally optimized parameters can cause negative
256 transfer, as local models may overfit to client-specific dis-
257 tortions. To assess cross-domain generalization, we intro-
258 duce a small server-held validation set \mathcal{D}_{val} that reflects di-
259 verse capture conditions across clients. Importantly, \mathcal{D}_{val}
260 is used solely for aggregation-time evaluation, without ac-
261 cessing raw client data.

262 Our objective is to learn a global model θ_{global} that (i)
263 achieves strong overall performance across heterogeneous
264 clients, (ii) remains robust under domain shift, and (iii) pre-
265 serves data privacy by sharing only model parameters.

266 3.2. Federated Learning

267 In standard federated learning, training proceeds in iterative
268 communication rounds between a central server and K dis-
269 tributed clients. At round t , the server broadcasts the current
270 global model parameters θ^t to participating clients. Each
271 client performs local optimization on its private dataset \mathcal{D}_k ,
272 producing updated parameters θ_k^{t+1} after one or more local
273 training epochs. The server then aggregates these locally
274 optimized parameters to obtain the next global model.

275 The most common aggregation rule, known as Federated
276 Averaging (FedAvg) [18], computes a data-weighted aver-
277 age of client updates:

$$278 \theta^{t+1} = \sum_{k=1}^K \frac{n_k}{\sum_{j=1}^K n_j} \theta_k^{t+1}.$$

279 This procedure approximates minimization of the global
280 empirical objective while preserving data privacy, since
281 only model parameters (or gradients) are communicated.

282 However, FedAvg implicitly assumes that client updates
283 are equally reliable outside their local domains. Under het-
284 erogeneous visual conditions, local updates may be biased
285 toward client-specific distortions, leading to suboptimal ag-
286 gregation and negative transfer [16]. In such non-IID set-
287 tings, uniform data-proportional weighting often fails to
288 capture the true cross-domain generalization capability of
289 individual clients. This limitation motivates the develop-
290 ment of aggregation strategies that explicitly incorporate
291 measures of client reliability during the update process.

3.3. ReSoFed

292 We introduce **ReSoFed**, a two-stage aggregation frame-
293 work that explicitly incorporates cross-domain reliability
294 into federated optimization. Rather than assuming uniform
295 contribution from all clients, ReSoFed evaluates the gen-
296 eralization behavior of locally trained client models using
297 a server-held validation set \mathcal{D}_{val} and integrates this signal
298 into the aggregation process. The framework consists of (i)
299 greedy soup-based client selection and (ii) reliability-aware
300 weighted aggregation, as illustrated in Figure 2.

301 **Greedy Soup-Based Client Selection:** Let $\{\theta_k\}_{k=1}^K$ denote
302 the locally optimized client models received at the server af-
303 ter training. We first estimate the reliability of each client
304 model using the validation set \mathcal{D}_{val} . Specifically, we com-
305 pute a reliability score:

$$307 r_k = \text{Eval}(\theta_k, \mathcal{D}_{val}),$$

308 where $\text{Eval}(\cdot)$ denotes a validation metric (e.g., accuracy in
309 our experiments). The reliability score r_k reflects how well
310 client k generalizes beyond its local data characteristics.

311 Motivated by the model soup paradigm [29], which
312 shows that weight-space averaging of independently fine-
313 tuned models can improve generalization without increas-
314 ing inference cost, we adapt greedy model souping to the
315 federated aggregation setting. Unlike conventional model
316 soups trained on centralized data, our approach performs
317 reliability-guided souping over distributed client models.

318 Clients are sorted in descending order of reliability. We
319 initialize the soup set S with the most reliable client:

$$320 S = \{k^*\}, \quad k^* = \arg \max_k r_k.$$

321 We then iteratively consider adding the remaining clients
322 to S . For a candidate client j , we compute the uniformly
323 averaged model over $S \cup \{j\}$:

$$324 \tilde{\theta}_{S \cup \{j\}} = \frac{1}{|S| + 1} \sum_{k \in S \cup \{j\}} \theta_k.$$

325 If the validation performance improves, i.e.,

$$326 \text{Eval}(\tilde{\theta}_{S \cup \{j\}}, \mathcal{D}_{val}) > \text{Eval}(\tilde{\theta}_S, \mathcal{D}_{val}),$$

327 then client j is added to S ; otherwise, it is discarded. This
328 greedy procedure continues until all clients have been eval-
329 uated.

330 This stage serves as a hard filtering mechanism, prevent-
331 ing poorly generalizing or incompatible client updates from
332 influencing the global model. By evaluating aggregated
333 models directly on \mathcal{D}_{val} , the selection process explicitly fa-
334 vors combinations that improve cross-domain robustness.

335 **Reliability-Aware Weighted Aggregation:** After identi-
336 fying the selected client subset S , we compute the final

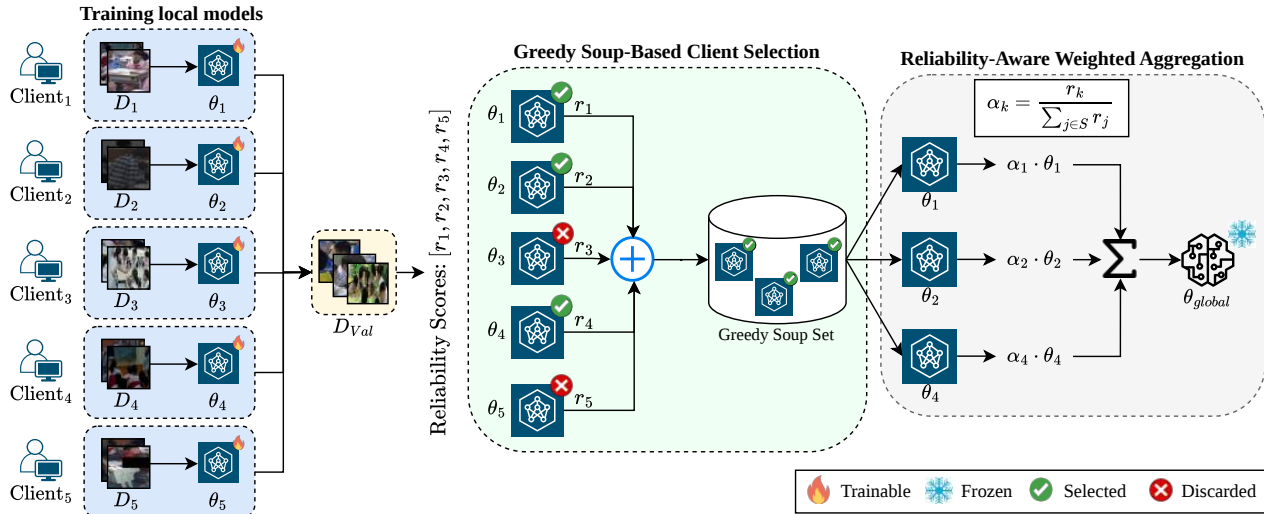


Figure 2. Overview of the proposed **ReSoFed** framework. Each client $k \in \{\text{Client}_1 \dots \text{Client}_5\}$ trains a local model θ_k on its private dataset \mathcal{D}_k under heterogeneous visual conditions (simulated via diverse augmentation regimes). Each locally trained model is evaluated on a diverse server-held validation set \mathcal{D}_{val} to compute a reliability score r_k . **Greedy Soup-Based Client Selection:** Clients are sorted in descending order of reliability, and a soup set S is initialized with the most reliable client $k^* = \arg \max_k r_k$. Additional clients are iteratively **included** only if the aggregated model improves validation performance. Poorly generalizing clients are **discarded**. **Reliability-Aware Weighted Aggregation:** The final global model is computed using reliability-proportional weights α_k resulting in $\theta_{global} = \sum_{k \in S} \alpha_k \theta_k$. Only model parameters are communicated, preserving data privacy while improving robustness under heterogeneous visual conditions.

337 global model using reliability-aware weighted aggregation
 338 over the selected clients. Instead of uniform averaging, we
 339 assign weights α_k proportional to the reliability scores:

$$340 \quad \alpha_k = \frac{r_k}{\sum_{j \in S} r_j}, \quad \text{for } k \in S.$$

341 The global model is then computed as:

$$342 \quad \theta_{global} = \sum_{k \in S} \alpha_k \theta_k.$$

343 This soft weighting mechanism further emphasizes
 344 clients that exhibit stronger cross-domain generalization
 345 while still allowing multiple reliable clients to contribute.
 346 Unlike standard FedAvg, which weights clients solely based
 347 on data size, ReSoFed directly incorporates validation-
 348 driven reliability into the aggregation process.

349 Importantly, the proposed framework does not modify
 350 local training objectives or optimization at the client level.
 351 All reliability estimation and selection steps are performed
 352 at the server using model parameters only, thereby preserv-
 353 ing privacy and introducing negligible computational over-
 354 head compared to standard federated aggregation. By combin-
 355 ing greedy soup-based client selection with reliability-
 356 weighted aggregation, ReSoFed effectively filters poorly
 357 generalizing clients while prioritizing models that demon-
 358 strate stronger cross-domain reliability. This reliability-
 359 guided aggregation strategy mitigates negative transfer un-

der heterogeneous visual conditions, leading to more stable
 global updates and improved overall performance.

4. Experiments 362

4.1. Setup 363

Dataset: We conduct experiments on the SCB (Student
 Classroom Behavior) dataset [30], which contains images
 collected from real classroom environments and annotated
 with bounding boxes representing student and teacher activ-
 ities. Using the provided annotations, we extract object-
 level image crops to construct an image classification
 dataset. The processed dataset contains 14 behavior cate-
 gories derived from classroom activities and scene ele-
 ments. The dataset exhibits notable class imbalance across
 behaviors, where frequent activities contain significantly
 more samples compared to rarer behaviors, resulting in a
 long-tailed distribution. In total, the resulting dataset con-
 tains over 90,000 labeled instances. It also captures diverse
 classroom interactions, viewpoints, and environmental con-
 ditions, providing a realistic benchmark for evaluating ro-
 bust classroom behavior recognition models. Representa-
 tive samples from the dataset are illustrated in Figure 3.

Federated Client Construction: To construct the feder-
 ated learning setting, the processed dataset is partitioned
 across five distributed clients, each representing an inde-
 pendent institution with its own local training data. Prior to
 partitioning, several low-frequency classes are removed due



Figure 3. Sample cropped instances extracted from the SCB dataset using bounding box annotations [30]. The dataset captures diverse classroom behaviors and visual conditions, including variations in viewpoint, occlusion, and lighting commonly observed in real educational environments.

to insufficient sample counts, ensuring that the remaining categories contain adequate instances for reliable training and evaluation. The filtered dataset is subsequently divided into training, validation, and tests subsets. The training portion is distributed across the five clients while maintaining balanced class representation, whereas the validation and test sets remain centralized at the server. This partitioning introduces client-specific data distributions, reflecting the non-IID characteristics commonly observed in federated learning scenarios. Under this configuration, each client performs local model optimization on its private dataset while sharing only model parameters with the central server, thereby preserving data privacy while enabling collaborative model training.

Heterogeneous Camera Simulation: To simulate heterogeneous visual capture conditions across educational institutions, each federated client is trained under a distinct augmentation regime designed to emulate common camera artifacts observed in real-world classroom recordings. These augmentations model variations frequently encountered in webcam-based data acquisition, including illumination changes, motion blur, compression artifacts, and partial occlusions arising from differences in camera hardware, lighting conditions, and recording environments. The specific augmentation configuration assigned to each client is summarized in Table 1. Such degradations are widely studied in robustness benchmarks and are known to affect computer vision models under distribution shift. In particular, corruption types such as brightness variation, blur, and compression artifacts are commonly used in robustness benchmarks such as ImageNet-C [8]. Similarly, occlusion-based perturbations have been extensively explored in vi-

Client	Simulated Condition	Reference
C_1	Clean reference capture	Standard baseline setting
C_2	Brightness shift	ImageNet-C [8]
C_3	Motion / defocus blur	ImageNet-C [8]
C_4	Compression artifacts	ImageNet-C [8]
C_5	Partial occlusion	Cutout [6], Random Erasing [31]

Table 1. Client-specific visual conditions used to simulate heterogeneous camera environments across federated institutions. Each client is trained under a distinct augmentation regime reflecting commonly observed real-world capture artifacts.

sual recognition through techniques such as Cutout [6] and Random Erasing [31]. In our setup, each client is assigned a specific augmentation configuration that approximates a distinct camera condition while preserving the semantic content of the images. This strategy introduces controlled cross-client visual heterogeneity, enabling systematic evaluation of federated learning methods under realistic domain shifts while maintaining consistent evaluation using a shared validation and test set.

Implementation Details: All experiments are conducted using five simulated federated clients, each trained under distinct visual augmentation settings to emulate heterogeneous capture conditions. Specifically, client heterogeneity is introduced through different augmentation regimes. Each client trains a local model using its augmented data partition without sharing raw images. We evaluate the proposed framework across multiple backbone architectures, including three convolutional neural networks (CNNs) [17] (*ResNet-50*, *EfficientNet-V2-S*, and *ConvNeXt-B*) and three transformer-based models [27] (*Swin-V2-T*, *Swin-V2-B*, and *Vision Transformer B-16*). All models are fine-tuned on the SCB dataset for the classroom action recognition task. Images are resized and normalized using standard ImageNet statistics. Local models are trained using the Adam optimizer with a learning rate of 1×10^{-4} . A small server-held validation set \mathcal{D}_{val} is used to estimate client reliability and guide the aggregation process in ReSoFed. This validation set is constructed to reflect diverse visual conditions and remains independent of all client training data. In our setup, reliability scores are computed using classification accuracy on \mathcal{D}_{val} . All experiments are implemented in *PyTorch* and executed on a single NVIDIA RTX 3090 GPU. We evaluate model performance using standard classification metrics, including Accuracy, Precision, Recall, and F1-score.

4.2. Results and Analysis

In this section, we present the evaluation results of the proposed framework, **ReSoFed**, on the SCB dataset for classroom action recognition. For comparison, we also report the performance of standard federated learning using Federated Averaging (FedAvg) under the same experimen-

Method	CNN Models					Transformer-based Models				
	Model	Acc (%)	Pre (%)	Rec (%)	F1 (%)	Model	Acc (%)	Pre (%)	Rec (%)	F1 (%)
Standard FL (FedAvg)	ResNet-50	73.82	77.34	73.67	72.83	ViT B-16	80.47	80.75	80.36	80.22
	EfficientNet-V2-S	78.93	80.31	78.79	78.45	Swin-V2-Tiny	81.23	82.56	81.10	80.80
	ConvNeXt-Base	81.82	81.97	80.78	81.37	Swin-V2-Base	81.96	83.26	81.73	81.41
ReSoFed	ResNet-50	78.50	78.38	78.44	78.06	ViT B-16	81.00	81.18	80.91	80.85
	EfficientNet-V2-S	81.61	81.49	81.57	81.22	Swin-V2-Tiny	81.68	83.01	81.58	81.18
	ConvNeXt-Base	82.75	83.00	82.68	82.35	Swin-V2-Base	82.35	82.55	82.29	82.16

Table 2. Performance comparison between standard federated learning (FedAvg) and the proposed ReSoFed framework on the SCB dataset. Results are reported across CNN and transformer architectures. **Bold** values indicate the best performance within each architecture group.

tal setting. Table 2 summarizes the classification performance across convolutional neural networks (CNNs) and transformer-based architectures to assess the effectiveness of ReSoFed across diverse model types.

The key observations from Table 2 are as follows: (1) ReSoFed consistently outperforms the standard federated aggregation (FedAvg) across all six models, including both CNN and transformer architectures, demonstrating the robustness and architecture-agnostic nature of the framework. (2) The performance improvements are particularly noticeable for CNN models. For instance, ResNet-50 shows a significant improvement in accuracy from 73.82% to 78.50%, while EfficientNet-V2-S improves from 78.93% to 81.61%, indicating that the proposed aggregation effectively mitigates negative transfer under heterogeneous visual conditions. (3) Although transformer architectures already achieve strong baseline performance under FedAvg, ReSoFed still provides consistent improvements, with Swin-V2-Base achieving the highest transformer performance at 82.35% accuracy. (4) Among all evaluated models, ConvNeXt-Base with ReSoFed achieves the best overall performance, reaching 82.75% accuracy and 82.35% F1-score, suggesting that modern CNN architectures benefit strongly from reliability-guided federated aggregation. (5) Overall, the consistent improvements across both architectures confirm that ReSoFed effectively handles client heterogeneity introduced by simulated camera conditions, validating the importance of reliability-guided client selection and weighted aggregation in federated learning settings.

4.3. Ablation Studies

To demonstrate the contribution of each component in the proposed ReSoFed framework, we conduct two ablation studies: (a) the effect of greedy soup-based client selection and (b) the effect of reliability-aware weighted aggregation. The ablations are performed on two representative backbone models, *ResNet-50* and *Swin-V2-Base*.

Effect of Greedy Soup-Based Client Selection: To evaluate the importance of the greedy client selection mechanism, we remove the greedy selection step and aggregate all client models using reliability-weighted averaging only.

Figure 4a presents the accuracy and F1 scores achieved by the two models with and without the greedy selection stage. The results show that removing the greedy selection procedure leads to noticeable performance degradation for both models. This indicates that greedy soup-based client selection plays a crucial role in filtering weak generalizing client models before aggregation, thereby reducing negative transfer and improving the robustness of the global model.

Effect of Reliability-Aware Weighted Aggregation: We next evaluate the impact of reliability-aware weighted aggregation. In this setting, we retain the greedy client selection stage but replace reliability-based weighting with uniform averaging among the selected clients. Figure 4b presents the performance before and after applying reliability weights during aggregation. The results demonstrate that reliability weighting provides additional performance gains beyond greedy selection alone. This confirms that reliability-aware weighting effectively prioritizes client models with stronger cross-domain generalization ability, further improving the performance of the aggregated model.

Overall, these ablations verify that both components of ReSoFed, greedy client selection and reliability-guided aggregation, contribute meaningfully to the final performance improvements observed in the proposed framework.

5. Discussion

The experimental results demonstrate that incorporating reliability-aware aggregation into federated learning can significantly improve model performance under heterogeneous visual conditions. By combining greedy model-soup-based client selection with reliability-weighted aggregation, ReSoFed mitigates the negative transfer effects that often arise when locally biased client models are aggregated indiscriminately. The consistent improvements observed across both convolutional and transformer-based backbones indicate that the proposed strategy is architecture-agnostic and can generalize across diverse model families.

Another important observation is that reliability estimation using a diverse server-held validation set provides an effective proxy for cross-domain generalization. Clients

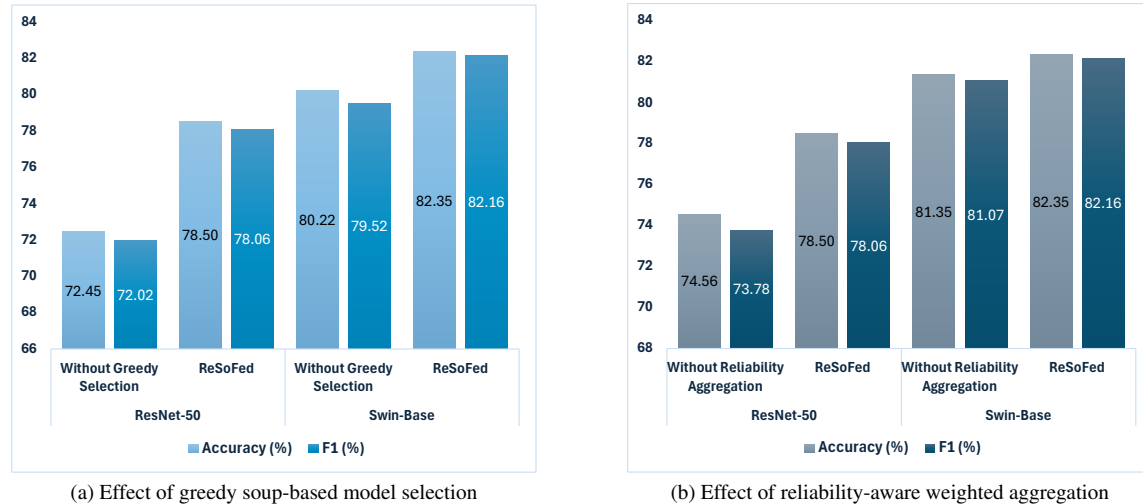


Figure 4. Ablation study of the ReSoFed framework. (a) Effect of greedy soup-based client selection, where removing greedy selection significantly degrades performance for both ResNet-50 and Swin-Base models. (b) Effect of reliability-aware weighted aggregation, showing that incorporating reliability scores during aggregation improves both accuracy and F1-score compared to uniform averaging.

537 whose locally trained models generalize well to this valida- 568
538 tion set are more likely to contribute positively to the global 569
539 model. This allows ReSoFed to selectively emphasize reli- 570
540 able clients while discarding harmful updates originating 571
541 from locally overfitted models. As a result, the global model 572
542 benefits from better overall classification performance. 573

543 Furthermore, the proposed approach operates entirely at 574
544 the aggregation stage without modifying the local training 575
545 procedure. This design choice makes ReSoFed easy to inte- 576
546 grate into existing federated learning pipelines. Since reli- 577
547 ability estimation and client selection occur only at the server 578
548 using model parameters, the approach preserves the core 579
549 privacy advantages of federated learning while introducing 580
550 minimal computational overhead. 581

551 **Limitations:** Despite its effectiveness, the proposed frame- 582
552 work has several limitations. First, the reliability estimation 583
553 depends on a server-held validation set representing diverse 584
554 visual conditions; if this set is not representative of the client 585
555 distributions, the reliability scores may not accurately cap- 586
556 ture cross-domain generalization. Second, our evaluation is 587
557 conducted in a simulated setting where heterogeneity is in- 588
558 troduced through controlled augmentations. Although these 589
559 augmentations are based on established robustness bench- 590
560 marks, real-world deployments may exhibit more complex 591
561 distribution shifts. Third, the greedy client selection stage 592
562 requires additional validation-time evaluation during aggre-
563 gation, which may increase computational cost as the num-
564 ber of clients grows.

565 6. Conclusion

566 In this work, we introduced **ReSoFed**, a reliability-guided 593
567 federated aggregation framework designed to improve ro- 594
595
596
597
598

599 bustness under heterogeneous visual conditions. The pro- 600
601 posed approach combines greedy model-soup-based client 602
603 selection with reliability-aware weighted aggregation, en- 604
605 abling the server to prioritize clients that demonstrate 606
607 stronger cross-domain generalization. By leveraging a 608
609 server-held validation set to estimate client reliability, Re- 610
611 SoFed mitigates the negative transfer that often arises when 612
613 locally biased models are aggregated indiscriminately. Ex- 614
615 perimental results on the SCB dataset for classroom action 616
617 recognition demonstrate that ReSoFed consistently outper- 618
619 forms standard federated learning, FedAvg, across multiple 620
621 backbone architectures, including both convolutional neural 622
623 networks and transformer-based models. These improve- 624
625 ments highlight the effectiveness of reliability-guided ag- 626
627 gregation for handling client heterogeneity while preserv- 628
629 ing the privacy of raw client data. 630

631 Future research directions include extending ReSoFed 632
633 to multi-round federated optimization settings, exploring 634
635 more advanced reliability estimation strategies beyond val- 636
637 idation accuracy, and evaluating the framework on larger 638
639 real-world federated settings with naturally occurring client 640
641 heterogeneity. Additionally, exploring alternative reliabil- 642
643 ity estimation strategies and more advanced client selection 644
645 mechanisms could further enhance the robustness of feder- 646
647 ated aggregation under diverse real-world data conditions. 648

649 References

- 650 [1] Kasra Borazjani, Payam Abdisarabshali, Naji Khosravan, 651
652 and Seyyedali Hosseinalipour. Redefining non-iid data in 653
654 federated learning for computer vision tasks: Migrating from 655
656 labels to embeddings for task-specific data distributions. 657
658 *arXiv preprint arXiv:2503.14553*, 2025. 1 659
660

- 599 [2] Mariah Bradford, Ibrahim Khebour, Nathaniel Blanchard,
600 and Nikhil Krishnaswamy. Automatic detection of collab-
601 orative states in small groups using multimodal features. In
602 *International Conference on Artificial Intelligence in Educa-*
603 *tion*, pages 767–773. Springer, 2023. 1
- 604 [3] H. Chen et al. Student behavior detection in classroom
605 videos based on improved yolov8. *Sensors*, 2023. 3
- 606 [4] Clayton Cohn, Caitlin Snyder, Joyce Horn Fonteles, Ashwin
607 TS, Justin Montenegro, and Gautam Biswas. A multimodal
608 approach to support teacher, researcher and ai collaboration
609 in stem+ c learning environments. *British Journal of Educa-*
610 *tional Technology*, 56(2):595–620, 2025. 1
- 611 [5] Wen Dai et al. Daisee: Towards user engagement recognition
612 in the wild. In *ACM International Conference on Multimodal*
613 *Interaction*, 2019. 3
- 614 [6] Terrance DeVries and Graham Taylor. Improved regulariza-
615 tion of convolutional neural networks with cutout. In *arXiv*
616 *preprint arXiv:1708.04552*, 2017. 6
- 617 [7] Dashan Gao, Xin Yao, and Qiang Yang. A survey
618 on heterogeneous federated learning. *arXiv preprint*
619 *arXiv:2210.04505*, 2022. 1, 2
- 620 [8] Dan Hendrycks and Thomas Dietterich. Benchmarking neu-
621 ral network robustness to common corruptions and perturba-
622 tions. In *ICLR*, 2019. 6
- 623 [9] Jiyue Huang, Chi Hong, Lydia Y. Chen, and Stefanie Roos.
624 Is shapley value fair? improving client selection for maver-
625 icks in federated learning. In *International Conference on*
626 *Machine Learning (ICML)*, 2021. 3
- 627 [10] Muhammad Rafsan Kabir, Rashidul Hassan Borshon, and
628 Riasat Khan. Federated learning for human activity recogni-
629 tion: Balancing privacy, efficiency, and accuracy through in-
630 novative aggregation techniques. *Array*, page 100462, 2025.
631 3
- 632 [11] Peter Kairouz and H Brendan McMahan. Advances and open
633 problems in federated learning. *Foundations and Trends in*
634 *Machine Learning*, 14(1-2):1–210, 2021. 1, 2
- 635 [12] S. N. Karimah et al. Automatic engagement estimation in
636 smart education: A review. *IEEE Access*, 2022. 3
- 637 [13] J. Li et al. Client selection strategies in federated learning:
638 A comprehensive survey. *Computer Networks*, 2024. 3
- 639 [14] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia
640 Smith. Federated learning: Challenges, methods, and future
641 directions. *IEEE Signal Processing Magazine*, 37(3):50–60,
642 2020. 1, 2
- 643 [15] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia
644 Smith. Federated optimization in heterogeneous networks.
645 *Proceedings of MLSys*, 2020. 3
- 646 [16] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and
647 Zhihua Zhang. On the convergence of fedavg on non-iid
648 data. In *International Conference on Learning Representa-*
649 *tions (ICLR)*, 2020. 4
- 650 [17] Zewen Li, Fan Liu, Wenjie Yang, Shouheng Peng, and Jun
651 Zhou. A survey of convolutional neural networks: analysis,
652 applications, and prospects. *IEEE Transactions on Neural*
653 *Networks and Learning Systems*, 33(12):6999–7019, 2021.
654 6
- [18] Brendan McMahan, Eider Moore, Daniel Ramage, Seth
Hampson, and Blaise Aguera y Arcas. Communication-
efficient learning of deep networks from decentralized data.
In *Proceedings of the International Conference on Artifi-*
cial Intelligence and Statistics (AISTATS), pages 1273–1282.
PMLR, 2017. 1, 2, 4
- [19] Mohammad Moshawrab et al. Reviewing federated learning
aggregation algorithms. *Electronics*, 12(10):2287, 2023. 3
- [20] Sindhuja Penchala et al. Learning in focus: Detecting behav-
ioral and collaborative engagement using vision transform-
ers. In *International Conference on Computer Vision in Ed-*
ucation, 2025. 3
- [21] Peng Qi et al. Model aggregation techniques in federated
learning: A systematic review. *Future Generation Computer*
Systems, 2024. 3
- [22] Liangqiong Qu, Yuyin Zhou, Paul Pu Liang, Yingda Xia,
Feifei Wang, Ehsan Adeli, Li Fei-Fei, and Daniel Rubin.
Rethinking architecture design for tackling data heterogene-
ity in federated learning. In *Proceedings of the IEEE/CVF*
Conference on Computer Vision and Pattern Recognition
(CVPR), pages 10061–10071, 2022. 1
- [23] Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary
Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and
Brendan McMahan. Adaptive federated optimization. In *In-*
ternational Conference on Learning Representations (ICLR),
2021. 3
- [24] Diganta Sengupta, Soumya Suvra Khan, Surajit Das, and
Debashis De. Fedel: Federated education learning for gen-
erating correlations between course outcomes and program
outcomes for internet of education things. *Internet of Things*,
25:101056, 2024. 1
- [25] Farhad M. Shiri et al. Detection of student engagement in
e-learning environments using hybrid deep learning models.
Journal of Artificial Intelligence, 2024. 3
- [26] B. Sun, Y. Wu, K. Zhao, et al. Student classroom behav-
ior dataset: A dataset for recognizing classroom behaviors.
Neural Computing and Applications, 2021. 3
- [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszko-
reit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia
Polosukhin. Attention is all you need. *Advances in Neural*
Information Processing Systems (NeurIPS), 30, 2017. 6
- [28] Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and
H. Vincent Poor. Tackling the objective inconsistency prob-
lem in heterogeneous federated optimization. In *Advances in*
Neural Information Processing Systems (NeurIPS), 2020. 3
- [29] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Re-
becca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos,
Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kor-
nblith, et al. Model soups: averaging weights of multi-
ple fine-tuned models improves accuracy without increas-
ing inference time. In *International Conference on Machine*
Learning (ICML), pages 23965–23998. PMLR, 2022. 4
- [30] Fan Yang. Scb-dataset: A dataset for detecting student class-
room behavior. *arXiv preprint arXiv:2304.02488*, 2023. 2,
5, 6
- [31] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and
Yi Yang. Random erasing data augmentation. In *AAAI*, 2020.
6