



HAL
open science

Decentralized Machine Learning with Centralized Performance Guarantees via Gibbs Algorithms

Yaiza Bermudez, Samir M. Perlaza, Iñaki Esnaola

► To cite this version:

Yaiza Bermudez, Samir M. Perlaza, Iñaki Esnaola. Decentralized Machine Learning with Centralized Performance Guarantees via Gibbs Algorithms. RR-9608, Inria. 2026. <hal-05471674>

HAL Id: hal-05471674

<https://inria.hal.science/hal-05471674v1>

Submitted on 28 Apr 2026

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY 4.0 - Attribution - International License



Decentralized Machine Learning with Centralized Performance Guarantees via Gibbs Algorithms

Yaiza Bermudez, Samir M. Perlaza, and Iñaki Esnaola

**RESEARCH
REPORT**

N° 9608

January 2026

Project-Team NEO

ISRN INRIA/RR--9608--FR+ENG

ISSN 0249-6399



Decentralized Machine Learning with Centralized Performance Guarantees via Gibbs Algorithms

Yaiza Bermudez, Samir M. Perlaza, and Iñaki Esnaola

Project-Team NEO

Research Report n° 9608 — January 2026 — 25 pages

Abstract: In this report, it is shown, for the first time, that centralized learning performance is achievable in decentralized learning without sharing the local datasets. Specifically, when clients adopt an empirical risk minimization with relative-entropy regularization (ERM-RER) learning framework and a forward-backward communication between clients is established, it suffices to share the locally obtained Gibbs measures to achieve the same performance as a centralized ERM-RER with access to all the datasets. The core mechanism is that the Gibbs measure produced by client k is used, as reference measure, by client $k + 1$. This effectively establishes a principled way to encode prior information through a reference measure. In particular, achieving centralized performance in the decentralized setting requires a specific scaling of the regularization factors with the local sample sizes. Overall, this result opens the door to novel decentralized learning paradigms that shift the collaboration strategy from sharing data to sharing the local inductive bias via the reference measures over the set of models.

Key-words: Decentralized machine learning; Gibbs algorithms; Gibbs measures; Empirical risk minimization; Relative-entropy regularization; KL divergence; Reference measure; Centralized performance guarantees; Peer-to-peer communication.

Yaiza Bermudez and Samir M. Perlaza are with INRIA, Centre Inria d'Université Côte d'Azur, Sophia Antipolis 06902, France. Samir M. Perlaza is also with the GAATI Laboratory at the Université de la Polynésie Française, Faaa 98702, French Polynesia. Iñaki Esnaola is with the School of Electrical and Electronic Engineering at The University of Sheffield, Sheffield S1 3JD, UK. Samir M. Perlaza and Iñaki Esnaola also are with the Electrical and Computer Engineering Department at Princeton University, Princeton N.J. 08544, USA. This research was supported in part by the European Commission through the H2020MSCA-RISE-2019 project 872172; the French National Agency for Research (ANR) through the Project ANR-21-CE25-0013 and the project ANR-22-PEFT-0010 of the France 2030 program PEPR Réseaux du Futur; and in part by the Agence de l'innovation de défense (AID) through the project UK-FR 2024352.

**RESEARCH CENTRE
SOPHIA ANTIPOLIS – MÉDITERRANÉE**

2004 route des Lucioles - BP 93
06902 Sophia Antipolis Cedex

Apprentissage décentralisé avec garanties de performance centralisée via des algorithmes de Gibbs

Résumé : Dans ce rapport, il est montré, pour la première fois, que des performances d'apprentissage centralisé sont atteignables en apprentissage décentralisé sans partage des jeux de données locaux. Plus précisément, lorsque les clients adoptent un cadre d'apprentissage de minimisation du risque empirique avec régularisation par entropie relative (ERM-RER) et qu'une communication aller-retour (forward-backward) est établie entre eux, il suffit de partager les mesures de Gibbs obtenues localement pour atteindre les mêmes performances qu'un ERM-RER centralisé ayant accès à l'ensemble des jeux de données. Le mécanisme clé est que la mesure de Gibbs produite par le client k est utilisée, comme mesure de référence, par le client $k+1$. Cela établit une manière principielle d'encoder une information a priori au moyen d'une mesure de référence. En particulier, l'obtention de performances centralisées dans le cadre décentralisé requiert une mise à l'échelle spécifique des facteurs de régularisation en fonction des tailles d'échantillon locales. Dans l'ensemble, ce résultat ouvre la voie à de nouveaux paradigmes d'apprentissage décentralisé qui déplacent la stratégie de collaboration : il ne s'agit plus de partager les données, mais de partager le biais inductif local via les mesures de référence sur l'ensemble des modèles.

Mots-clés : Apprentissage décentralisé; Algorithmes de Gibbs ; Mesures de Gibbs ; Minimisation du risque empirique ; Régularisation par entropie relative ; Divergence de Kullback–Leibler ; Mesure de référence ; Garanties de performance centralisée ; Communication pair-à-pair.

Contents

1	Introduction	4
1.1	Contributions	4
1.2	Notation and Preliminaries	5
2	Supervised Machine Learning	6
3	Gibbs Algorithms	7
4	Main Result	8
4.1	Communication Protocol	8
4.2	Decentralized Algorithms	9
4.3	Centralized performance guarantees	10
5	Conclusion and Final Remarks	13
	References	14
	Appendices	17
A	Preliminaries for the Proof of Theorem 4	17
B	Proof of Theorem 4	20
C	Proof of Theorem 6	22
D	Complementary Results	23

1 Introduction

Decentralized learning studies how a collection of clients can collaboratively tune a learning algorithm by communicating only over a network, without explicitly exchanging raw datasets. This setting extends early work on distributed and asynchronous optimization, where coordination is achieved through local computations and intermittent message passing [1, 2]. It becomes particularly relevant when a central coordinator is unavailable or undesirable, or when data transfers are impractical due to bandwidth, latency, ownership, privacy, or regulatory constraints [3]. A standard benchmark for collaborative learning is the centralized regime in which all local datasets are pooled and a single training procedure is run on the aggregated data. While conceptually simple, the pooled-data benchmark is often unachievable in decentralized environments due to communication constraints and/or restricted disclosure of local datasets [3, 4]. This benchmark is revisited through the lens of *Gibbs algorithms*, i.e., probability measures on the model space. Gibbs measures arise naturally as solutions of empirical risk minimizations with relative-entropy regularization (ERM-RER) [5–8]. This viewpoint also connects to exponential-weights predictors and PAC-Bayesian posteriors, which reason directly in terms of distributions on hypotheses [9–14]. Beyond their variational interpretation, Gibbs measures also capture the long-run distribution of stochastic gradient methods under suitable regimes [15–18]. From this standpoint, a complementary line of work studies Gibbs measures as solutions to ERM-RER problems and its extensions [6–8]. Other studies focus on change-of-measure techniques to quantify the variation of an expectation when the underlying probability measure changes [5, 19]. These developments provide tools to interpret and manipulate Gibbs measures as first-class objects in learning systems, and to reason about how information is transported through probability measures rather than through datasets.

1.1 Contributions

This report shows that centralized performance guarantees can be achieved in a decentralized system through a strategic design of (i) the *reference measures* and (ii) the *regularization factors* that define the clients’ Gibbs algorithm. More precisely, a peer-to-peer communication protocol is introduced, in which each client transmits its Gibbs probability measure to its successor that adopts it as reference measure. This mechanism embeds information from datasets into the learning process without explicitly transmitting such datasets. A closed-form expression is obtained for the resulting decentralized Gibbs probability measures, together with conditions under which it coincides with the Gibbs measure induced by the centralized pooled-data benchmark.

The report is organized as follows. Section 2 presents the supervised learning setting in a decentralized system and fixes the problem formulation. Section 3 defines Gibbs conditional probability measures and their interpretation within the context of ERM-RER. Section 4 presents the communication protocol and the main results establishing centralized-performance guarantees. Finally, Section 5

concludes and discusses practical considerations, including the impact of finite-rate communication distortions on the transmission of probability measures. Detailed proofs are provided in the appendices.

1.2 Notation and Preliminaries

The set of all probability measures on the measurable space $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$ is denoted by $\Delta(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$, or simply $\Delta(\mathcal{X})$. Moreover, the set of probability measures in $\Delta(\mathcal{X})$ that are absolutely continuous with respect to Q is denoted by $\Delta_Q(\mathcal{X})$. Using this notation, a conditional probability measure is defined hereunder.

Definition 1 (Conditional Probability). *A family $P_{Y|X} \triangleq (P_{Y|X=x})_{x \in \mathcal{X}}$ of elements of $\Delta(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$ indexed by \mathcal{X} is said to be a conditional probability measure if, for all sets $\mathcal{B} \in \mathcal{F}_{\mathcal{Y}}$, the map*

$$\begin{cases} \mathcal{X} \longrightarrow [0, 1] \\ x \longmapsto P_{Y|X=x}(\mathcal{B}) \end{cases} \quad (1)$$

is Borel measurable. The set of such conditional probability measures is denoted by $\Delta(\mathcal{Y}|\mathcal{X})$.

Given two measures P and Q on the same measurable space, the notation $P \ll Q$ stands for “the measure P is absolutely continuous with respect to Q ”. Using this notation, the relative entropy or KL-divergence is defined hereunder.

Definition 2 (Relative Entropy). *Given a probability measure P and a σ -finite measure Q , both on the same measurable space, with $P \ll Q$. The relative entropy of P with respect to Q is*

$$D(P \parallel Q) \triangleq \int \frac{dP}{dQ}(x) \log \left(\frac{dP}{dQ}(x) \right) dQ(x). \quad (2)$$

Lemma 1. *Consider three probability measures P, Q , and R on the same measurable space such that $R \ll P \ll Q$, it then follows that*

$$\int \log \left(\frac{dP}{dQ}(\theta) \right) dR(\theta) = D(R \parallel Q) - D(R \parallel P). \quad (3)$$

Proof: Under the absolute continuity assumptions, i.e., $R \ll P \ll Q$, from Definition 2 it follows that

$$\int \log \left(\frac{dP}{dQ}(\theta) \right) dR(\theta) = \int \log \left(\frac{dP}{dR}(\theta) \frac{dR}{dQ}(\theta) \right) dR(\theta) \quad (4)$$

$$= \int -\log \left(\frac{dR}{dP}(\theta) \right) + \log \left(\frac{dR}{dQ}(\theta) \right) dR(\theta) \quad (5)$$

$$= D(R \parallel Q) - D(R \parallel P), \quad (6)$$

where the equality in (4) follows from [20, Theorem 5] and the equality in (5) follows from [20, Theorem 4]. \blacksquare

2 Supervised Machine Learning

Consider a decentralized learning system in which K clients collaboratively tune their local learning algorithms by communicating with each other. For all $k \in \{1, 2, \dots, K\}$, let \mathcal{M}_k , \mathcal{X}_k and \mathcal{Y}_k , with $\mathcal{M}_k \subseteq \mathbb{R}^{d_k}$ and $d_k \in \mathbb{N}$, be sets of *models*, *patterns*, and *labels*, respectively, at client k . The training data available for client k consists of n_k data points $(x_{k,1}, y_{k,1}), (x_{k,2}, y_{k,2}), \dots, (x_{k,n_k}, y_{k,n_k})$, which are elements of the set $\mathcal{Z}_k \triangleq \mathcal{X}_k \times \mathcal{Y}_k$. Such data points form the local training dataset, denoted by $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$, which can be explicitly written as

$$\mathbf{z}_k \triangleq ((x_{k,1}, y_{k,1}), (x_{k,2}, y_{k,2}), \dots, (x_{k,n_k}, y_{k,n_k})). \quad (7)$$

The dataset obtained by the aggregation of all local datasets, denoted by \mathbf{z}_0 , satisfies

$$\mathbf{z}_0 \triangleq (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K) \in \mathcal{Z}_1^{n_1} \times \mathcal{Z}_2^{n_2} \times \dots \times \mathcal{Z}_K^{n_K} \quad (8)$$

$$= ((x_{0,1}, y_{0,1}), (x_{0,2}, y_{0,2}), \dots, (x_{0,n_0}, y_{0,n_0})). \quad (9)$$

Hence, the total number of data points, denoted by $n_0 \in \mathbb{N}$, satisfies

$$n_0 \triangleq \sum_{k=1}^K n_k. \quad (10)$$

Given a model $\theta \in \mathcal{M}_k$ for client k , the loss induced by such a model with respect to a data point $(x, y) \in \mathcal{Z}_k$ is $\ell_k(x, y, \theta)$, where the function

$$\ell_k : \mathcal{Z}_k \times \mathcal{M}_k \rightarrow [0, +\infty), \quad (11)$$

is referred to as the *loss function* of client k . Such a loss function is measurable with respect to the measurable spaces $(\mathcal{Z}_k \times \mathcal{M}_k, \mathcal{F}_k)$ and $([0, +\infty), \mathcal{B}([0, +\infty)))$, where \mathcal{F}_k is a given σ -field on $\mathcal{Z}_k \times \mathcal{M}_k$ and $\mathcal{B}([0, +\infty))$ is the Borel σ -field on $[0, +\infty)$. The *empirical risk* induced by such a model $\theta \in \mathcal{M}_k$, with respect to the dataset \mathbf{z}_k in (7), is determined by the function

$$\mathbb{L}_k : \begin{cases} \mathcal{Z}_k^{n_k} \times \mathcal{M}_k \longrightarrow [0, +\infty) \\ (\mathbf{z}_k, \theta) \longmapsto \frac{1}{n_k} \sum_{i=1}^{n_k} \ell_k(x_{k,i}, y_{k,i}, \theta), \end{cases} \quad (12)$$

where the function ℓ_k is defined in (11). In the following, a supervised machine learning algorithm is represented by a conditional probability measure, as defined hereunder.

Definition 3 (Algorithm). *For all $k \in \{1, 2, \dots, K\}$, a conditional probability measure $P_{\theta_k | \mathbf{z}_k} \in \Delta(\mathcal{M}_k | (\mathcal{X}_k \times \mathcal{Y}_k)^{n_k})$ is said to represent a supervised machine learning algorithm.*

Let $P_{\theta_k | \mathbf{z}_k} \in \Delta(\mathcal{M}_k | (\mathcal{X}_k \times \mathcal{Y}_k)^{n_k})$ be an algorithm. Hence, the instance of such an algorithm trained upon the dataset \mathbf{z}_k in (7) is denoted by $P_{\theta_k | \mathbf{z}_k = \mathbf{z}_k}$, which is simply a probability measure in $\Delta(\mathcal{M}_k)$.

A class of algorithms that are central in this work are known as Gibbs algorithms. These algorithms are introduced in the following section.

3 Gibbs Algorithms

Gibbs algorithms are represented by Gibbs conditional probability measures. Such conditional probability measures are parametrized by the empirical risk function L_k ; a σ -finite measure $Q_k \in \Delta(\mathcal{M}_k)$; and a dataset $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$, with $k \in \{1, 2, \dots, K\}$. In order to define such Gibbs conditional measures, consider the following function:

$$\mathsf{K}_{k, Q_k, \mathbf{z}_k} : \begin{cases} \mathbb{R} \longrightarrow \mathbb{R} \\ t \longmapsto \log \left(\int \exp(t L_k(\mathbf{z}_k, \boldsymbol{\theta}_k)) dQ_k(\boldsymbol{\theta}_k) \right). \end{cases} \quad (13)$$

Under the assumption that the reference measure Q_k is a probability measure, the function $\mathsf{K}_{k, Q_k, \mathbf{z}_k}$ in (13) is the cumulant generating function of the random variable $L_k(\mathbf{z}_k, \boldsymbol{\theta}_k)$, for some fixed dataset $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$, when the model $\boldsymbol{\theta}_k$ is sampled from Q_k . Using this notation, the definition of Gibbs conditional probability measures is presented hereunder.

Definition 4. *Given the function L_k in (12); a σ -finite measure Q_k ; and a $\lambda_k \in (0, +\infty)$, with $k \in \{1, 2, \dots, K\}$, the probability measure $P_{\boldsymbol{\theta}_k | \mathbf{Z}_k}^{(Q_k, \lambda_k)} \in \Delta(\mathcal{M}_k | \mathcal{Z}_k^{n_k})$ is said to be an (L_k, Q_k, λ_k) -Gibbs conditional probability measure if*

$$\forall \mathbf{z}_k \in \mathcal{S}_k, \mathsf{K}_{k, Q_k, \mathbf{z}_k} \left(\frac{-1}{\lambda_k} \right) < +\infty; \quad (14)$$

for some set $\mathcal{S}_k \subseteq \mathcal{Z}_k^{n_k}$; and for all $(\mathbf{z}_k, \boldsymbol{\theta}_k) \in \mathcal{S}_k \times \text{supp } Q_k$,

$$\frac{dP_{\boldsymbol{\theta}_k | \mathbf{Z}_k}^{(Q_k, \lambda_k)}}{dQ_k}(\boldsymbol{\theta}_k) = \exp \left(\frac{-1}{\lambda_k} L_k(\mathbf{z}_k, \boldsymbol{\theta}_k) - \mathsf{K}_{k, Q_k, \mathbf{z}_k} \left(\frac{-1}{\lambda_k} \right) \right), \quad (15)$$

where the function $\mathsf{K}_{k, Q_k, \mathbf{z}_k}$ is defined in (13).

In the following, if the probability measure from which datasets are sampled by client k is $P_{\mathbf{Z}_k} \in \Delta(\mathcal{Z}_k^{n_k})$, the set \mathcal{S}_k in Definition 4 is chosen to be the support of $P_{\mathbf{Z}_k}$. Note that, while $P_{\boldsymbol{\theta}_k | \mathbf{Z}_k}^{(Q_k, \lambda_k)}$ in (15) is referred to as a Gibbs conditional probability measure, the measure $P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_k, \lambda_k)}$, obtained by conditioning upon a given dataset $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$, is referred to as a Gibbs probability measure. Consider the functional $\mathsf{R}_{k, \mathbf{z}_k}$ defined as follows,

$$\mathsf{R}_{k, \mathbf{z}_k} : \begin{cases} \Delta(\mathcal{M}_k) \longrightarrow [0, +\infty) \\ P \longmapsto \int L_k(\mathbf{z}_k, \boldsymbol{\theta}) dP(\boldsymbol{\theta}), \end{cases} \quad (16)$$

where the function L_k is defined in (12). The Gibbs probability measure $P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_k, \lambda_k)}$ in (15) can be related to the following optimization problem:

$$\min_{P \in \Delta_{Q_k}(\mathcal{M}_k)} \mathsf{R}_{k, \mathbf{z}_k}(P) + \lambda_k D(P \| Q_k). \quad (17)$$

The following lemma formalizes this connection.

Lemma 2. *Assume that the optimization problem in (17) admits a solution. Then the probability measure $P_{\Theta_k | \mathbf{Z}_k}^{(Q_k, \lambda_k)}$ in (15) is the unique solution.*

Proof: The proof follows from [19, Lemma 1]. ■

This result has also been reported for other f -divergences in [7, 8]. Interestingly, the conditional probability measure $P_{\Theta_k | \mathbf{Z}_k}^{(Q_k, \lambda_k)}$ in (15) is the long-run distribution of a stochastic gradient descent algorithm [18]. In statistical learning, such a distribution is often referred to as the *Gibbs algorithm* [21].

Another optimization problem that is closely related to the probability measure $P_{\Theta_k | \mathbf{Z}_k}^{(Q_k, \lambda_k)}$ in (15) is the following:

$$\min_{P \in \Delta_{Q_k}(\mathcal{M}_k)} R_{k, \mathbf{z}_k}(P) \tag{18a}$$

$$\text{s.t. } D(P \parallel Q_k) \leq \gamma_k, \tag{18b}$$

for some $\gamma_k > 0$. The following lemma establishes the connection.

Lemma 3. *Assume that the optimization problem in (18) admits a solution and that λ_k is such that*

$$D\left(P_{\Theta_k | \mathbf{Z}_k}^{(Q_k, \lambda_k)} \parallel Q_k\right) = \gamma_k. \tag{19}$$

Then, the probability measure $P_{\Theta_k | \mathbf{Z}_k}^{(Q_k, \lambda_k)}$ in (15) is the unique solution.

Proof: The proof follows from [19, Lemma 4]. ■

Lemma 3 implies that the probability measure $P_{\Theta_k | \mathbf{Z}_k}^{(Q_k, \lambda_k)}$ in (15) minimizes the training empirical risk over all probability measures in the following neighborhood of Q_k ,

$$\{P \in \Delta_{Q_k}(\mathcal{M}_k) : D(P \parallel Q_k) \leq \gamma_k\}. \tag{20}$$

4 Main Result

The main result of this work (Theorem 4) is presented in Subsection 4.2. In order to present such a result, the peer-to-peer communication protocol used by the clients is introduced. The section ends by stating the necessary conditions under which centralized performance is obtained.

4.1 Communication Protocol

Figure 1 depicts the forward–backward peer-to-peer communication protocol used in this work. In the forward direction (blue arrows), for all $k \in \{1, 2, \dots, K-1\}$, client k transmits to client $k+1$ the (L_k, Q_k, λ_k) -Gibbs probability measure $P_{\Theta_k | \mathbf{Z}_k}^{(Q_k, \lambda_k)}$ in (15), obtained as the solution to the ERM-RER problem in (17). Client $k+1$ adopts this transmitted measure as its reference measure Q_{k+1} in (17). This choice induces a nested structure of the reference measure, as

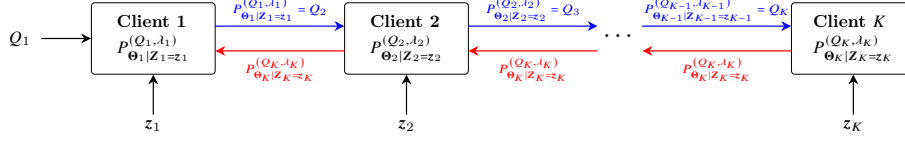


Figure 1: Nested Structure: the Gibbs measure produced by client k becomes the reference measure Q_{k+1} used by client $k + 1$.

$Q_{k+1} = P_{\theta_k | Z_k = z_k}^{(Q_k, \lambda_k)}$, with $k \in \{1, 2, \dots, K-1\}$. This is formalized later in Theorem 4. The backward direction (red arrows) disseminates the final Gibbs probability measure. More specifically, once client K has computed its (L_K, Q_K, λ_K) -Gibbs probability measure $P_{\theta_K | Z_K = z_K}^{(Q_K, \lambda_K)}$ in (15), this measure is transmitted back along the chain, from client K to client $K - 1$, then from client $K - 1$ to client $K - 2$, and so on until client 1. This backward transmission provides all clients with access to the same final Gibbs algorithm.

4.2 Decentralized Algorithms

The main result of this work is presented in terms of the Gibbs algorithms obtained through the communication protocol described in Section 4.1. In particular, for all $k \in \{1, 2, \dots, K\}$, in such a protocol, the (L_k, Q_k, λ_k) -Gibbs probability measure

$$P_{\theta_k | Z_k = z_k}^{(Q_k, \lambda_k)} \in \Delta(\mathcal{M}_k), \quad (21a)$$

with $\lambda_k > 0$, represents the Gibbs algorithm of client k trained upon its local dataset \mathbf{z}_k defined in (7). The corresponding reference measure Q_k is specified recursively as

$$Q_k = \begin{cases} Q_1, & \text{if } k = 1, \\ P_{\theta_{k-1} | Z_{k-1} = z_{k-1}}^{(Q_{k-1}, \lambda_{k-1})}, & \text{if } k \geq 2, \end{cases} \quad (21b)$$

From Definition 4, for all $\theta \in \text{supp } Q_1$, it follows that

$$\frac{dP_{\theta_k | Z_k = z_k}^{(Q_k, \lambda_k)}(\theta)}{dQ_k}(\theta) = \frac{\exp\left(-\frac{1}{\lambda_k} L_k(z_k, \theta)\right)}{\int \exp\left(-\frac{1}{\lambda_k} L_k(z_k, \nu)\right) dQ_k(\nu)}. \quad (21c)$$

From Lemma 3, for all $k \in \{1, 2, \dots, K\}$, it follows that such an algorithm is optimal in the sense that the corresponding probability measure $P_{\theta_k | Z_k = z_k}^{(Q_k, \lambda_k)}$ in (21) minimizes the expected training empirical risk with respect to \mathbf{z}_k , over all probability measures in the following neighborhood of Q_k ,

$$\mathcal{G}_k \triangleq \left\{ P \in \Delta_{Q_1}(\mathcal{M}_k) : D(P \parallel Q_k) \leq D\left(P_{\theta_k | Z_k = z_k}^{(Q_k, \lambda_k)} \parallel Q_k\right) \right\}. \quad (22)$$

The main result of this work is presented by the following theorem.

Theorem 4. For all $k \in \{1, 2, \dots, K\}$, consider the (L_k, Q_k, λ_k) -Gibbs probability measure $P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_k, \lambda_k)}$ in (21). Then, for all $\theta \in \text{supp } Q_1$,

$$\frac{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_k, \lambda_k)}(\theta)}{dQ_1} = \frac{\exp\left(\sum_{j=1}^k \frac{-1}{\lambda_j} L_j(\mathbf{z}_j, \theta)\right)}{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \nu)\right) dQ_1(\nu)}. \quad (23)$$

Proof: The proof is presented in Appendix B. ■

The choice of reference measures Q_1, Q_2, \dots, Q_K in (21b) induces a *nested* structure. Under this structure, the training performed by client k uses only its local dataset, while the influence of the previous clients datasets is carried out through the reference measure Q_k . The relevance of this nested structure in which client k shares its Gibbs probability measure (algorithm) with its $k + 1$ neighbor, is made clear by [19, Lemma 1]. More specifically, the conditional measure $P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_k, \lambda_k)}$ in (23) is the unique solution to (17) and, simultaneously, the unique minimizer of the following optimization problem:

$$\min_{P \in \Delta_{Q_1}(\mathcal{M}_k)} \int \left(\sum_{j=1}^k \frac{1}{\lambda_j} L_j(\mathbf{z}_j, \theta) \right) dP(\theta) + D(P \| Q_1). \quad (24)$$

The following corollary of Theorem 4 formalizes this observation.

Corollary 5. Under the assumption that the measures Q_1, Q_2, \dots, Q_K satisfy (21b), the solutions to the optimization problems in (17) and (24) are unique and coincide with the (L_k, Q_k, λ_k) -Gibbs probability measure $P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_k, \lambda_k)}$ in (23).

Given $k > 1$, the optimization problem in (17) depends exclusively on the local training dataset \mathbf{z}_k . While the reference Q_k in (21b) depends on the training datasets of the $k - 1$ previous clients, the dependence is not direct. Solving (17), requires access to the probability measure $Q_k \in \Delta(\mathcal{M}_k)$, but not access to the training datasets of all previous clients. This is because for fixed training datasets, Q_k is simply a probability measure on the model space. In contrast, the optimization problem in (24) depends on the training datasets of client j , for all $j \in \{1, 2, \dots, k\}$, while the reference measure Q_1 does not depend on any training dataset. The fact that problems (17) and (24) share the same solution unveils an important observation: providing client k with a reference measure Q_k of the form in (21b) reproduces the effect of having access to the training datasets of the $k - 1$ previous clients, without transmitting those datasets.

4.3 Centralized performance guarantees

An important observation is that a strategic choice of $\lambda_1, \lambda_2, \dots, \lambda_K$ in (23) can lead to achieve the same Gibbs probability distribution as in a setting in which the training datasets of all clients are available at all clients. This describes a decentralized system whose distributed nature does not limit achieving the same Gibbs algorithm that one would have obtained if all the training data were available at all clients. The following theorem formalizes this observation.

Theorem 6. Assume that the loss functions in (11) satisfy $\ell_1 = \ell_2 = \dots = \ell_K = \ell$ and for all $k \in \{1, \dots, K\}$,

$$\lambda_k = \frac{n_0}{n_k} \lambda_0, \quad (25)$$

for some $\lambda_0 > 0$ and some loss function ℓ . Consider some measures Q_1, Q_2, \dots, Q_K satisfying (21b). Then, for all $\theta \in \text{supp } Q_1$, it follows that,

$$\frac{dP_{\Theta_K | Z_K = z_K}^{(Q_K, \lambda_K)}}{dP_{\Theta_K | Z_0 = z_0}^{(Q_1, \lambda_0)}}(\theta) = 1, \quad (26)$$

where the probability measure $P_{\Theta_K | Z_K = z_K}^{(Q_K, \lambda_K)}$ is defined in (23); the measure $P_{\Theta_K | Z_0 = z_0}^{(Q_1, \lambda_0)}$ satisfies for all $\theta \in \text{supp } Q_1$,

$$\frac{dP_{\Theta_K | Z_0 = z_0}^{(Q_1, \lambda_0)}}{dQ_1}(\theta) = \frac{\exp\left(\frac{-1}{n_0 \lambda_0} \sum_{i=1}^{n_0} \ell(x_{0,i}, y_{0,i}, \theta)\right)}{\int \exp\left(\frac{-1}{n_0 \lambda_0} \sum_{i=1}^{n_0} \ell(x_{0,i}, y_{0,i}, \nu)\right) dQ_1(\nu)}; \quad (27)$$

and $(x_{0,i}, y_{0,i})$ are data points of the aggregated dataset z_0 in (8).

Proof: The proof is presented in Appendix C. ■

The relevance of Theorem 6 is highlighted by the following observations. Under the assumptions of Theorem 6, in particular that the loss functions satisfy $\ell_1 = \ell_2 = \dots = \ell_K = \ell$, the equality in (25) together with [21, Lemma 4] allows rewriting the optimization problem in (24) as

$$\min_{P \in \Delta_{Q_1}(\mathcal{M}_K)} \int \frac{1}{n_0} \sum_{i=1}^{n_0} \ell(x_{0,i}, y_{0,i}, \theta) dP(\theta) + \lambda_0 D(P \parallel Q_1), \quad (28)$$

which requires access to the training datasets of all clients. Interestingly, from [19, Lemma 1], it follows that the probability measure $P_{\Theta_K | Z_0 = z_0}^{(Q_1, \lambda_0)}$ in (27) is the solution to (28). More importantly, from [19, Lemma 4], if λ_0 is chosen such that

$$D\left(P_{\Theta_K | Z_0 = z_0}^{(Q_1, \lambda_0)} \parallel Q_1\right) = \gamma_0, \quad (29)$$

for some $\gamma_0 > 0$, the probability measure $P_{\Theta_K | Z_0 = z_0}^{(Q_1, \lambda_0)}$ in (27) is also the solution to the following optimization problem:

$$\min_{P \in \Delta_{Q_1}(\mathcal{M}_K)} \int \frac{1}{n_0} \sum_{i=1}^{n_0} \ell(x_{0,i}, y_{0,i}, \theta) dP(\theta) \quad (30a)$$

$$\text{s.t. } D(P \parallel Q_1) \leq \gamma_0. \quad (30b)$$

Hence, given that from (26), the measures $P_{\Theta_K | Z_K = z_K}^{(Q_K, \lambda_K)}$ and $P_{\Theta_K | Z_0 = z_0}^{(Q_1, \lambda_0)}$ are identical, the nested structure induced by the choice of reference measures in (21b) achieves the same Gibbs probability measure as a centralized system in which all training

datasets are available at client K . More specifically, under the forward–backward communication protocol in Section 4.1, after $K - 1$ forward messages (blue arrows in Figure 1) client K obtains the Gibbs algorithm that minimizes the empirical risk with respect to all training datasets within the following neighborhood of Q_1 ,

$$\mathcal{G}_0 \triangleq \left\{ P \in \Delta_{Q_1}(\mathcal{M}_K) : D(P \parallel Q_1) \leq D\left(P_{\Theta_K|Z_0=z_0}^{(Q_1, \lambda_K)} \parallel Q_1\right) \right\}. \quad (31)$$

The backward dissemination (red arrows in Figure 1) then provides each client such a Gibbs algorithm.

Figure 2 provides a representation of the nested structure induced by the choice of reference measure in (21) under the assumption in (25): it depicts the successive displacement of the reference measure along the forward pass of the communication protocol, together with the corresponding sets \mathcal{G}_k in (22), with $k \in \{1, 2, \dots, K\}$, induced by the relative-entropy regularization. Starting from the reference measure Q_1 , client 1 outputs the Gibbs probability measure $P_{\Theta_1|Z_1=z_1}^{(Q_1, \lambda_1)}$ in (27), represented by the point Q_2 , which belongs to the set \mathcal{G}_1 . This locally obtained measure then becomes the reference measure for client 2, whose Gibbs probability measure $P_{\Theta_2|Z_2=z_2}^{(Q_2, \lambda_2)}$ in (27) is represented by Q_3 and lies in \mathcal{G}_2 . Iterating this construction yields the sequence $Q_1 \rightarrow Q_2 \rightarrow Q_3 \rightarrow Q_4 \rightarrow Q_5$, where each point Q_{k+1} represents $P_{\Theta_k|Z_k=z_k}^{(Q_k, \lambda_k)}$ in (27) and each dashed region \mathcal{G}_k describes the admissible neighborhood around Q_k enforced by (22).

In particular, the figure shows the sets \mathcal{G}_0 in (31) and \mathcal{G}_k in (22) together with the probability measures $P_{\Theta_k|Z_k=z_k}^{(Q_k, \lambda_k)}$ in (23), with $k \in \{1, 2, \dots, K\}$, and $P_{\Theta_K|Z_0=z_0}^{(Q_1, \lambda_K)}$ in (27). The endpoint Q_5 is annotated as $Q_5 = P_{\Theta_K|Z_0=z_0}^{(Q_1, \lambda_K)}$, highlighting that the nested reference-measure construction produced by the communication protocol reaches the same Gibbs probability measure as the direct construction on the aggregated dataset within the set \mathcal{G}_0 in (31).

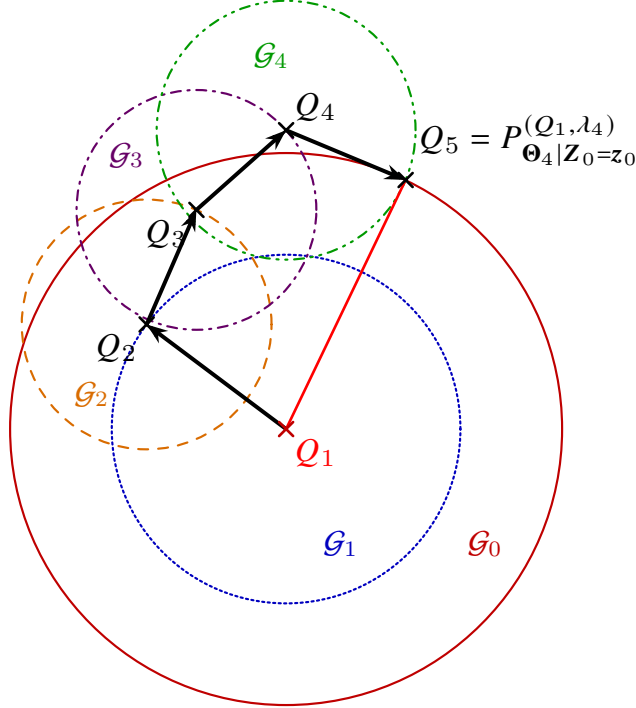


Figure 2: Geometric representation, under the assumptions in (25), of the sets \mathcal{G}_0 in (31) and \mathcal{G}_k in (22), with $k \in \{1, 2, \dots, 4\}$. The points Q_{k+1} represent the Gibbs probability measures $P_{\Theta_k | Z_k = z_k}^{(Q_k, \lambda_k)}$ in (23), and the endpoint is $Q_5 = P_{\Theta_4 | Z_0 = z_0}^{(Q_1, \lambda_4)}$ as in (27).

5 Conclusion and Final Remarks

In this work, it has been shown that there exists a choice of reference measures (Q_1, Q_2, \dots, Q_K) and regularization factors $(\lambda_1, \lambda_2, \dots, \lambda_K)$ in a decentralized machine learning scenario such that the same performance of a centralized system in which all training datasets are aggregated and made available is achievable. The construction of such regularization factors is rather simple. The regularization factor of client k shall be the product of a strictly positive real (common to all clients) and the ratio of data points in the training dataset of client k and the total number of data points across the training datasets of all clients. The reference measure of client k , $k > 1$, is the Gibbs measure (Gibbs algorithm) used by client $k - 1$ to sample its models. The first client uses a given reference Q_1 . Client K can back propagate its achieved Gibbs probability measure (Gibbs algorithm) across all clients in such a way that all clients achieve the same performance. That is, they sample their models from the same probability measure, which minimizes the empirical risk with respect to the training datasets of all clients within a neighborhood of Q_1 .

This choice of reference measures leads to a nested structure whose construction requires the transmission of $K - 1$ probability measures with common support [6, Lemma 3]. The practical construction of this nested structure faces several challenges. The most prominent is a consequence of data transmission within a finite number of bits and delay constraints. This implies that the Gibbs probability measure transmitted by client k might be received by client $k + 1$ subject to distortion. The impact of such a distortion in the construction of the nested structure is an open problem and has not been taken into account in this work. Another challenge is related to the fact that the Gibbs measures to be transmitted might have supports that are significantly large and thus, the amount of data to be transmitted might be comparable to transmitting the training datasets. Nonetheless, the transmission of a probability measure is more privacy preserving than the actual transmission of training datasets. In a nutshell, several important theoretical discoveries have been reported in this work concerning decentralized learning.

References

- [1] J. N. Tsitsiklis, D. P. Bertsekas, and M. Athans, “Distributed asynchronous deterministic and stochastic gradient optimization algorithms,” *IEEE Transactions on Automatic Control*, vol. 31, no. 9, pp. 803–812, Sept. 1986.
- [2] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*, 1st ed. Englewood Cliffs, NJ: Prentice-Hall, 1989.
- [3] P. Kairouz, B. H. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, R. G. L. d’Oliveira, S. E. Rouayheb, D. Evans, J. Gardner, Z. Garrett, A. Gascón, B. Ghazi, P. B. Gibbons, M. Gruteser, Z. Harchaoui, C. He, L. He, Z. Huo, B. Hutchinson, J. Hsu, M. Jaggi, T. Javidi, G. Joshi, M. Khodak, J. Konečný, A. Korolova, F. Koushanfar, S. Koyejo, T. Lepoint, Y. Liu, P. Mittal, M. Mohri, R. Nock, A. Özgür, R. Pagh, M. Raykova, H. Qi, D. Ramage, R. Raskar, D. Song, W. Song, S. U. Stich, Z. Sun, A. T. Suresh, F. Tramèr, P. Vepakomma, J. Wang, L. Xiong, Z. Xu, Q. Yang, F. X. Yu, H. Yu, and S. Zhao, *Advances and Open Problems in Federated Learning*. Now Publishers, 2021, vol. 14, no. 1–2.
- [4] C. Dwork, “Differential privacy,” in *Proceedings of the 33rd International Colloquium on Automata, Languages and Programming (ICALP)*, vol. 4052, Venice, Italy, Jul. 2006, pp. 1–12.
- [5] I. Csiszár, “ I -divergence geometry of probability distributions and minimization problems,” *The Annals of Probability*, vol. 3, no. 1, pp. 146–158, Feb. 1975.

- [6] S. M. Perlaza, G. Bisson, I. Esnaola, A. Jean-Marie, and S. Rini, “Empirical risk minimization with relative entropy regularization,” *IEEE Transactions on Information Theory*, vol. 70, no. 7, pp. 5122 – 5161, Jul. 2024.
- [7] F. Daunas, I. Esnaola, S. M. Perlaza, and H. V. Poor, “Equivalence of empirical risk minimization to regularization on the family of f -divergences,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, Athens, Greece, Jul. 2024, pp. 759–764.
- [8] —, “Asymmetry of the relative entropy in the regularization of empirical risk minimization,” *IEEE Transactions on Information Theory*, vol. 71, no. 8, pp. 6198–6226, Aug. 2025.
- [9] N. Cesa-Bianchi and G. Lugosi, *Prediction, Learning, and Games*, 1st ed. New York, NY, USA: Cambridge University Press, 2006.
- [10] D. A. McAllester, “Some PAC-Bayesian theorems,” *Machine Learning*, vol. 37, no. 3, pp. 355–363, Dec. 1999.
- [11] M. Seeger, “PAC-Bayesian generalisation error bounds for Gaussian process classification,” *Journal of Machine Learning Research*, vol. 3, pp. 233–269, Oct. 2002.
- [12] J. Langford and J. Shawe-Taylor, “PAC-Bayes and margins,” in *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 15, Vancouver, Canada, Dec. 2002, pp. 439–446.
- [13] O. Catoni, *PAC-Bayesian Supervised Classification: The Thermodynamics of Statistical Learning*, 1st ed. Beachwood, OH, USA: Institute of Mathematical Statistics Lecture Notes - Monograph Series, 2007, vol. 56.
- [14] P. Alquier, “User-friendly introduction to PAC-Bayes bounds,” *Foundations and Trends in Machine Learning*, vol. 17, no. 2, pp. 174–303, 2024.
- [15] M. Welling and Y. W. Teh, “Bayesian learning via stochastic gradient Langevin dynamics,” in *Proceedings of the 28th International Conference on Machine Learning (ICML)*, Bellevue, Washington, USA, Jun. 2011, pp. 681–688.
- [16] S. Mandt, M. D. Hoffman, and D. M. Blei, “Stochastic gradient descent as approximate Bayesian inference,” *Journal of Machine Learning Research*, vol. 18, no. 1, p. 4873–4907, Jan. 2017.
- [17] M. Raginsky, A. Rakhlin, and M. Telgarsky, “Non-convex learning via stochastic gradient Langevin dynamics: A nonasymptotic analysis,” in *Proceedings of the Conference on Learning Theory (COLT)*, vol. 65, Amsterdam, Netherlands, Jul. 2017, pp. 1674–1703.

- [18] W. Azizian, F. Lutzeler, J. Malick, and P. Mertikopoulos, “What is the long-run distribution of stochastic gradient descent? A large deviations analysis,” in *Proceedings of the International Conference on Machine Learning (ICML)*, Vienna, Austria, July 2024, pp. 2168 – 2229.
- [19] S. M. Perlaza and G. Bisson, “Variations on the expectation due to changes in the probability measure,” *Entropy*, vol. 27, no. 8:865, pp. 1–20, Aug. 2025.
- [20] Y. Bermudez, G. Bisson, I. Esnaola, and S. M. Perlaza, “Proofs for folklore theorems on the Radon-Nikodym derivative,” INRIA, Centre Inria d’Université Côte d’Azur, Sophia Antipolis, France, Tech. Rep. RR-9591, July 2025.
- [21] S. M. Perlaza, I. Esnaola, G. Bisson, and H. V. Poor, “On the validation of Gibbs algorithms: Training datasets, test datasets and their aggregation,” in *Proceedings of the International Symposium on Information Theory (ISIT)*, Taipei, Taiwan, June 2023, pp. 328–333.
- [22] H. Jeffreys, “An invariant form for the prior probability in estimation problems,” *Proceedings of the Royal Society of London. Series A, Mathematical and Physical Sciences*, vol. 186, no. 1007, pp. 453–461, Sept. 1946.

Appendices

A Preliminaries for the Proof of Theorem 4

The following lemma establishes a relationship between the $(\mathbf{L}_k, \mathcal{Q}_k, \lambda_k)$ -Gibbs probability measure in (21) and an $(\mathbf{L}_k, \mathcal{Q}_1, \lambda_k)$ -Gibbs probability measure $P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_1, \lambda_k)} \in \Delta(\mathcal{M}_k)$.

Lemma 7. *For all $k \in \{1, 2, \dots, K\}$, consider the $(\mathbf{L}_k, \mathcal{Q}_k, \lambda_k)$ -Gibbs probability measure $P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_k, \lambda_k)}$ in (21). Then, for all $\boldsymbol{\theta} \in \text{supp } \mathcal{Q}_1$,*

$$\frac{dP_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_k, \lambda_k)}(\boldsymbol{\theta})}{d\mathcal{Q}_k} = \frac{dP_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_1, \lambda_k)}(\boldsymbol{\theta}) \exp(C_k)}{d\mathcal{Q}_1}, \quad (32)$$

where the measure $P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_1, \lambda_k)}$ is an $(\mathbf{L}_k, \mathcal{Q}_1, \lambda_k)$ -Gibbs probability measure and $C_k \in \mathbb{R}$ satisfies

$$C_k = \log \left(\frac{\int \exp \left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} \mathbf{L}_i(\mathbf{z}_i, \tilde{\mathbf{v}}) \right) d\mathcal{Q}_1(\tilde{\mathbf{v}}) \left(\exp \left(\mathbf{K}_{k, \mathcal{Q}_1, \mathbf{z}_k} \left(\frac{-1}{\lambda_k} \right) \right) \right)}{\int \exp \left(\sum_{j=1}^k \frac{-1}{\lambda_j} \mathbf{L}_j(\mathbf{z}_j, \mathbf{v}) \right) d\mathcal{Q}_1(\mathbf{v})} \right), \quad (33)$$

where the functional $\mathbf{K}_{k, \mathcal{Q}_1, \mathbf{z}_k}$ is defined as in (13).

Proof: The proof of (32) is done by induction on k .

Base case ($k = 2$): From Definition 4 and \mathcal{Q}_2 in (21b), for all $\boldsymbol{\theta} \in \text{supp } \mathcal{Q}_1$, it follows that

$$\begin{aligned} & \frac{dP_{\boldsymbol{\theta}_2 | \mathbf{Z}_2 = \mathbf{z}_2}^{(P_{\boldsymbol{\theta}_1 | \mathbf{Z}_1 = \mathbf{z}_1}^{(\mathcal{Q}_1, \lambda_1)}, \lambda_2)}(\boldsymbol{\theta})}{dP_{\boldsymbol{\theta}_1 | \mathbf{Z}_1 = \mathbf{z}_1}^{(\mathcal{Q}_1, \lambda_1)}} \\ &= \exp \left(\frac{-1}{\lambda_2} \mathbf{L}_2(\mathbf{z}_2, \boldsymbol{\theta}) - \mathbf{K}_{2, P_{\boldsymbol{\theta}_1 | \mathbf{Z}_1 = \mathbf{z}_1}^{(\mathcal{Q}_1, \lambda_1)}, \mathbf{z}_2} \left(\frac{-1}{\lambda_2} \right) \right) \end{aligned} \quad (34)$$

$$= \frac{\exp \left(\frac{-1}{\lambda_2} \mathbf{L}_2(\mathbf{z}_2, \boldsymbol{\theta}) \right)}{\int \exp \left(\frac{-1}{\lambda_2} \mathbf{L}_2(\mathbf{z}_2, \mathbf{v}) \right) \frac{dP_{\boldsymbol{\theta}_1 | \mathbf{Z}_1 = \mathbf{z}_1}^{(\mathcal{Q}_1, \lambda_1)}(\mathbf{v})}{d\mathcal{Q}_1}(\mathbf{v}) d\mathcal{Q}_1(\mathbf{v})} \quad (35)$$

$$= \frac{\exp \left(\frac{-1}{\lambda_2} \mathbf{L}_2(\mathbf{z}_2, \boldsymbol{\theta}) \right)}{\int \exp \left(\frac{-1}{\lambda_2} \mathbf{L}_2(\mathbf{z}_2, \mathbf{v}) \right) \left(\frac{\exp \left(\frac{-1}{\lambda_1} \mathbf{L}_1(\mathbf{z}_1, \mathbf{v}) \right)}{\int \exp \left(\frac{-1}{\lambda_1} \mathbf{L}_1(\mathbf{z}_1, \boldsymbol{\theta}') \right) d\mathcal{Q}_1(\boldsymbol{\theta}')} \right) d\mathcal{Q}_1(\mathbf{v})} \quad (36)$$

$$= \frac{\exp \left(\frac{-1}{\lambda_2} \mathbf{L}_2(\mathbf{z}_2, \boldsymbol{\theta}) \right) \int \exp \left(\frac{-1}{\lambda_1} \mathbf{L}_1(\mathbf{z}_1, \boldsymbol{\theta}') \right) d\mathcal{Q}_1(\boldsymbol{\theta}')}{\int \exp \left(\frac{-1}{\lambda_2} \mathbf{L}_2(\mathbf{z}_2, \mathbf{v}) \right) \exp \left(\frac{-1}{\lambda_1} \mathbf{L}_1(\mathbf{z}_1, \mathbf{v}) \right) d\mathcal{Q}_1(\mathbf{v})} \quad (37)$$

$$= \frac{\exp\left(\frac{-1}{\lambda_2}L_2(z_2, \theta)\right) \int \exp\left(\frac{-1}{\lambda_1}L_1(z_1, \theta')\right) dQ_1(\theta') \exp\left(K_{2, Q_1, z_2}\left(\frac{-1}{\lambda_2}\right)\right)}{\int \exp\left(\frac{-1}{\lambda_2}L_2(z_2, \nu)\right) \exp\left(\frac{-1}{\lambda_1}L_1(z_1, \nu)\right) dQ_1(\nu) \exp\left(K_{2, Q_1, z_2}\left(\frac{-1}{\lambda_2}\right)\right)} \quad (38)$$

$$= \frac{dP_{\Theta_2|Z_2=z_2}^{(Q_1, \lambda_2)}(\theta)}{dQ_1}(\theta) \frac{\exp\left(K_{1, Q_1, z_1}\left(\frac{-1}{\lambda_1}\right) + K_{2, Q_1, z_2}\left(\frac{-1}{\lambda_2}\right)\right)}{\int \exp\left(\frac{-1}{\lambda_1}L_1(z_1, \nu) + \frac{-1}{\lambda_2}L_2(z_2, \nu)\right) dQ_1(\nu)}, \quad (39)$$

where (35) follows from [20, Theorem 2]; and (36) follows from Definition 4. The equality in (39) corresponds to (32) in the case of $k = 2$.

Induction hypothesis: Assume that (32) holds for $k \geq 2$.

Induction step ($k \rightarrow k+1$): It remains to prove (32) for $k+1$. Consider client $k+1$ and its reference measure Q_k in (21b). From Definition 4, it holds that for all $\theta \in \text{supp } Q_1$,

$$\frac{dP_{\Theta_{k+1}|Z_{k+1}}^{(Q_{k+1}, \lambda_{k+1})}(\theta)}{dQ_{k+1}}(\theta) = \frac{\exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(z_{k+1}, \theta)\right)}{\int \exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(z_{k+1}, \nu)\right) dQ_{k+1}(\nu)} \quad (40)$$

$$= \frac{\exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(z_{k+1}, \theta)\right)}{\int \exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(z_{k+1}, \nu)\right) \frac{dQ_{k+1}(\nu)}{dQ_k(\nu)} dQ_k(\nu)} \quad (41)$$

$$= \frac{\exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(z_{k+1}, \theta)\right)}{\int \exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(z_{k+1}, \nu)\right) \frac{dP_{\Theta_k|Z_k=z_k}^{(Q_k, \lambda_k)}(\nu)}{dQ_k(\nu)} dQ_k(\nu)}, \quad (42)$$

where the equality in (41) follows from [20, Theorem 2]; and the equality in (42) follows from (21b). By the induction hypothesis, it holds that, for all $\nu \in \text{supp } Q_1$

$$\begin{aligned} & \frac{dP_{\Theta_k|Z_k=z_k}^{(Q_k, \lambda_k)}(\nu)}{dQ_k}(\nu) \\ &= \frac{dP_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}(\nu)}{dQ_1}(\nu) \frac{\int \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i}L_i(z_i, \tilde{\nu})\right) dQ_1(\tilde{\nu}) \left(\exp\left(K_{k, Q_1, z_k}\left(\frac{-1}{\lambda_k}\right)\right)\right)}{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i}L_i(z_i, \tilde{\theta})\right) dQ_1(\tilde{\theta})} \end{aligned} \quad (43)$$

$$= \frac{\exp\left(\frac{-1}{\lambda_k}L_k(z_k, \nu)\right) \int \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i}L_i(z_i, \tilde{\nu})\right) dQ_1(\tilde{\nu})}{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i}L_i(z_i, \tilde{\theta})\right) dQ_1(\tilde{\theta})}. \quad (44)$$

Substituting (44) into the denominator of (42) yields for all $\theta \in \text{supp } Q_1$,

$$\frac{dP_{\Theta_{k+1}|Z_{k+1}}^{(Q_{k+1}, \lambda_{k+1})}(\theta)}{dQ_{k+1}}(\theta)$$

$$\begin{aligned}
 &= \frac{\exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(\mathbf{z}_{k+1}, \boldsymbol{\theta})\right)}{\int \exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(\mathbf{z}_{k+1}, \boldsymbol{\nu})\right) \exp\left(\frac{-1}{\lambda_k}L_k(\mathbf{z}_k, \boldsymbol{\nu})\right) dQ_k(\boldsymbol{\nu})} \\
 &\quad \cdot \frac{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\theta}})\right) dQ_1(\tilde{\boldsymbol{\theta}})}{\int \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\nu}})\right) dQ_1(\tilde{\boldsymbol{\nu}})} \tag{45}
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{\exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(\mathbf{z}_{k+1}, \boldsymbol{\theta})\right)}{\int \exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(\mathbf{z}_{k+1}, \boldsymbol{\nu}) + \frac{-1}{\lambda_k}L_k(\mathbf{z}_k, \boldsymbol{\nu})\right) \frac{dQ_k}{dQ_{k-1}}(\boldsymbol{\nu}) \dots \frac{dQ_3}{dQ_2}(\boldsymbol{\nu}) \frac{dQ_2}{dQ_1}(\boldsymbol{\nu}) dQ_1(\boldsymbol{\nu})} \\
 &\quad \cdot \frac{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\theta}})\right) dQ_1(\tilde{\boldsymbol{\theta}})}{\int \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\nu}})\right) dQ_1(\tilde{\boldsymbol{\nu}})}, \tag{46}
 \end{aligned}$$

where (46) follows from [20, Theorem 4]. Moreover, for all $\boldsymbol{\nu} \in \text{supp } Q_1$, it follows that,

$$\prod_{j=2}^k \left(\frac{dQ_j}{dQ_{j-1}}(\boldsymbol{\nu}) \right) \tag{47}$$

$$= \prod_{j=2}^k \left(\frac{dP_{\boldsymbol{\theta}_{j-1} | \mathbf{Z}_{j-1} = \mathbf{z}_{j-1}}^{(Q_{j-1}, \lambda_{j-1})}}{dQ_{j-1}}(\boldsymbol{\nu}) \right) \tag{48}$$

$$= \prod_{j=1}^{k-1} \left(\frac{dP_{\boldsymbol{\theta}_j | \mathbf{Z}_j = \mathbf{z}_j}^{(Q_j, \lambda_j)}}{dQ_j}(\boldsymbol{\nu}) \right) \tag{49}$$

$$= \prod_{j=1}^{k-1} \left(\frac{\exp\left(\frac{-1}{\lambda_j}L_j(\mathbf{z}_j, \boldsymbol{\nu})\right) \int \exp\left(\sum_{i=1}^{j-1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\nu}})\right) dQ_1(\tilde{\boldsymbol{\nu}})}{\int \exp\left(\sum_{i=1}^j \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\theta}})\right) dQ_1(\tilde{\boldsymbol{\theta}})} \right) \tag{50}$$

$$= \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \boldsymbol{\nu})\right) \prod_{j=1}^{k-1} \left(\frac{\int \exp\left(\sum_{i=1}^{j-1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\nu}})\right) dQ_1(\tilde{\boldsymbol{\nu}})}{\int \exp\left(\sum_{i=1}^j \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\theta}})\right) dQ_1(\tilde{\boldsymbol{\theta}})} \right) \tag{51}$$

$$= \frac{\exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \boldsymbol{\nu})\right)}{\int \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\theta}})\right) dQ_1(\tilde{\boldsymbol{\theta}})}, \tag{52}$$

where the equality in (49) follows from (21b); the equality in (49) follows by changing the indices in (47); and the equality in (50) follows from the induction hypothesis.

Substituting (52) into (46) yields for all $\boldsymbol{\theta} \in \text{supp } Q_1$,

$$\frac{dP_{\boldsymbol{\theta}_{k+1} | \mathbf{Z} = \mathbf{z}_{k+1}}^{(Q_{k+1}, \lambda_{k+1})}}{dQ_{k+1}}(\boldsymbol{\theta})$$

$$\begin{aligned}
&= \frac{\exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(\mathbf{z}_{k+1}, \boldsymbol{\theta})\right)}{\int \exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(\mathbf{z}_{k+1}, \boldsymbol{\nu}) + \frac{-1}{\lambda_k}L_k(\mathbf{z}_k, \boldsymbol{\nu})\right) \left(\exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \boldsymbol{\nu})\right)\right) dQ_1(\boldsymbol{\nu})} \\
&\cdot \left(\int \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\nu}})\right) dQ_1(\tilde{\boldsymbol{\nu}})\right) \frac{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\theta}})\right) dQ_1(\tilde{\boldsymbol{\theta}})}{\int \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\nu}})\right) dQ_1(\tilde{\boldsymbol{\nu}})} \quad (53)
\end{aligned}$$

$$\begin{aligned}
&= \frac{\exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(\mathbf{z}_{k+1}, \boldsymbol{\theta})\right) \int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\theta}})\right) dQ_1(\tilde{\boldsymbol{\theta}})}{\int \exp\left(\sum_{i=1}^{k+1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \boldsymbol{\nu})\right) dQ_1(\boldsymbol{\nu})} \quad (54)
\end{aligned}$$

$$\begin{aligned}
&= \frac{\exp\left(\frac{-1}{\lambda_{k+1}}L_{k+1}(\mathbf{z}_{k+1}, \boldsymbol{\theta})\right) \int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\theta}})\right) dQ_1(\tilde{\boldsymbol{\theta}})}{\int \exp\left(\sum_{i=1}^{k+1} \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \boldsymbol{\nu})\right) dQ_1(\boldsymbol{\nu})} \\
&\frac{\exp\left(\mathbb{K}_{k+1, Q_1, \mathbf{z}_{k+1}}\left(\frac{-1}{\lambda_{k+1}}\right)\right)}{\exp\left(\mathbb{K}_{k+1, Q_1, \mathbf{z}_{k+1}}\left(\frac{-1}{\lambda_{k+1}}\right)\right)} \quad (55)
\end{aligned}$$

$$\begin{aligned}
&= \frac{dP_{\boldsymbol{\theta}_{k+1} | \mathbf{Z}_{k+1} = \mathbf{z}_{k+1}}^{(Q_1, \lambda_{k+1})}(\boldsymbol{\theta})}{dQ_1} \\
&\frac{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\theta}})\right) dQ_1(\tilde{\boldsymbol{\theta}}) \left(\exp\left(\mathbb{K}_{k+1, Q_1, \mathbf{z}_{k+1}}\left(\frac{-1}{\lambda_{k+1}}\right)\right)\right)}{\int \exp\left(\sum_{j=1}^{k+1} \frac{-1}{\lambda_j}L_j(\mathbf{z}_j, \boldsymbol{\nu})\right) dQ_1(\boldsymbol{\nu})}, \quad (56)
\end{aligned}$$

where (54) follows by simplifying the ratio of products in the second factor in (53).

The second factor in (56) does not depend on $\boldsymbol{\theta}$ and is denoted by

$$C_{k+1} \triangleq \log \left(\frac{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i}L_i(\mathbf{z}_i, \tilde{\boldsymbol{\theta}})\right) dQ_1(\tilde{\boldsymbol{\theta}}) \left(\exp\left(\mathbb{K}_{k+1, Q_1, \mathbf{z}_{k+1}}\left(\frac{-1}{\lambda_{k+1}}\right)\right)\right)}{\int \exp\left(\sum_{j=1}^{k+1} \frac{-1}{\lambda_j}L_j(\mathbf{z}_j, \boldsymbol{\nu})\right) dQ_1(\boldsymbol{\nu})} \right). \quad (57)$$

Combining (56) and (57) yields for all $\boldsymbol{\theta} \in \text{supp } Q_1$,

$$\frac{dP_{\boldsymbol{\theta}_{k+1} | \mathbf{Z} = \mathbf{z}_{k+1}}^{(Q_{k+1}, \lambda_{k+1})}(\boldsymbol{\theta})}{dQ_{k+1}} = \frac{dP_{\boldsymbol{\theta}_{k+1} | \mathbf{Z}_{k+1} = \mathbf{z}_{k+1}}^{(Q_1, \lambda_{k+1})}(\boldsymbol{\theta})}{dQ_1} \exp(C_{k+1}), \quad (58)$$

which corresponds to (32) in the case of $k+1$ and concludes the proof. \blacksquare

B Proof of Theorem 4

From Lemma 7, for all $\boldsymbol{\theta} \in \text{supp } Q_1$, it follows that,

$$\frac{dP_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_k, \lambda_k)}(\boldsymbol{\theta})}{dQ_k} = \frac{dP_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_1, \lambda_k)}(\boldsymbol{\theta})}{dQ_1} \exp(C_k), \quad (59)$$

with

$$C_k = \log \left(\frac{\int \exp \left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \tilde{\mathbf{v}}) \right) dQ_1(\tilde{\mathbf{v}}) \left(\exp \left(K_{k, Q_1, \mathbf{z}_k} \left(\frac{-1}{\lambda_k} \right) \right) \right)}{\int \exp \left(\sum_{j=1}^k \frac{-1}{\lambda_j} L_j(\mathbf{z}_j, \mathbf{v}) \right) dQ_1(\mathbf{v})} \right). \quad (60)$$

From [6, Lemma 3] and [20, Theorem 2] the Radon–Nikodym derivatives $\frac{dQ_1}{dQ_k}$ and $\frac{dQ_k}{dQ_1}$ are well defined. From the chain rule of the Radon–Nikodym derivative [20, Theorem 4], for all $\boldsymbol{\theta} \in \text{supp } Q_1$, it follows that,

$$\frac{dP_{\boldsymbol{\Theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_k, \lambda_k)}(\boldsymbol{\theta})}{dQ_k} = \frac{dP_{\boldsymbol{\Theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_1, \lambda_k)}(\boldsymbol{\theta})}{dQ_1} \frac{dQ_1}{dQ_k}(\boldsymbol{\theta}). \quad (61)$$

Moreover, from the multiplicative inverse of the Radon–Nikodym derivative [20, Theorem 5], by combining (59) and (61), for all $\boldsymbol{\theta} \in \text{supp } Q_1$, it follows that,

$$\begin{aligned} & \frac{dP_{\boldsymbol{\Theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_k, \lambda_k)}(\boldsymbol{\theta})}{dQ_1} \\ &= \frac{dQ_k}{dQ_1}(\boldsymbol{\theta}) \frac{dP_{\boldsymbol{\Theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_1, \lambda_k)}(\boldsymbol{\theta})}{dQ_1} \exp(C_k) \end{aligned} \quad (62)$$

$$= \frac{dQ_k}{dQ_{k-1}}(\boldsymbol{\theta}) \frac{dQ_{k-1}}{dQ_{k-2}}(\boldsymbol{\theta}) \dots \frac{dQ_3}{dQ_2}(\boldsymbol{\theta}) \frac{dQ_2}{dQ_1}(\boldsymbol{\theta}) \frac{dP_{\boldsymbol{\Theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_1, \lambda_k)}(\boldsymbol{\theta})}{dQ_1} \exp(C_k) \quad (63)$$

$$= \prod_{j=2}^k \left(\frac{dQ_j}{dQ_{j-1}}(\boldsymbol{\theta}) \right) \frac{dP_{\boldsymbol{\Theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_1, \lambda_k)}(\boldsymbol{\theta})}{dQ_1} \exp(C_k), \quad (64)$$

where (63) follows from [20, Theorem 4]. From Lemma 7 and (52), for all $\boldsymbol{\theta} \in \text{supp } Q_1$ it follows that

$$\prod_{j=2}^k \left(\frac{dQ_j}{dQ_{j-1}}(\boldsymbol{\theta}) \right) = \frac{\exp \left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \boldsymbol{\theta}) \right)}{\int \exp \left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \mathbf{v}) \right) dQ_1(\mathbf{v})}. \quad (65)$$

Plugging (65) into (64) yield

$$\begin{aligned} & \frac{dP_{\boldsymbol{\Theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_k, \lambda_k)}(\boldsymbol{\theta})}{dQ_1} \\ &= \frac{\exp \left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \boldsymbol{\theta}) \right)}{\int \exp \left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \mathbf{v}) \right) dQ_1(\mathbf{v})} \frac{dP_{\boldsymbol{\Theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_1, \lambda_k)}(\boldsymbol{\theta})}{dQ_1} \exp(C_k) \\ &= \frac{\exp \left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \boldsymbol{\theta}) \right)}{\int \exp \left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \mathbf{v}) \right) dQ_1(\mathbf{v})} \frac{dP_{\boldsymbol{\Theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_1, \lambda_k)}(\boldsymbol{\theta})}{dQ_1} \end{aligned} \quad (66)$$

$$\frac{\int \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \tilde{\mathbf{v}})\right) dQ_1(\tilde{\mathbf{v}}) \left(\exp\left(K_{k, Q_1, \mathbf{z}_k}\left(\frac{-1}{\lambda_k}\right)\right)\right)}{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \mathbf{v})\right) dQ_1(\mathbf{v})} \quad (67)$$

$$= \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \boldsymbol{\theta})\right) \frac{dP_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_1, \lambda_k)}(\boldsymbol{\theta})}{dQ_1} \frac{\exp\left(K_{k, Q_1, \mathbf{z}_k}\left(\frac{-1}{\lambda_k}\right)\right)}{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \mathbf{v})\right) dQ_1(\mathbf{v})} \quad (68)$$

$$= \frac{\exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \boldsymbol{\theta})\right) \exp\left(\frac{-1}{\lambda_k} L_k(\mathbf{z}_k, \boldsymbol{\theta})\right)}{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \mathbf{v})\right) dQ_1(\mathbf{v})} \quad (69)$$

$$= \frac{\exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \boldsymbol{\theta})\right)}{\int \exp\left(\sum_{i=1}^k \frac{-1}{\lambda_i} L_i(\mathbf{z}_i, \mathbf{v})\right) dQ_1(\mathbf{v})}, \quad (70)$$

where (67) follows from (60); and (69) follows from Definition 4. This completes the proof.

C Proof of Theorem 6

For all $k \in \{1, \dots, K\}$, under the assumption in (25), it holds that $\lambda_k = \frac{n_0}{n_k} \lambda_0$. From Theorem 4, it follows that for all $\boldsymbol{\theta} \in \text{supp } Q_1$,

$$\frac{dP_{\boldsymbol{\theta}_K | \mathbf{Z}_K = \mathbf{z}_K}^{(Q_K, \lambda_K)}(\boldsymbol{\theta})}{dQ_1} = \frac{\exp\left(\sum_{k=1}^K \frac{-n_k}{n_0 \lambda_0} L_k(\mathbf{z}_k, \boldsymbol{\theta})\right)}{\int \exp\left(\sum_{k=1}^K \frac{-n_k}{n_0 \lambda_0} L_k(\mathbf{z}_k, \mathbf{v})\right) dQ_1(\mathbf{v})} \quad (71)$$

$$= \frac{\exp\left(\sum_{k=1}^K \frac{-n_k}{n_0 \lambda_0} \frac{1}{n_k} \sum_{i=1}^{n_k} \ell_k(x_{k,i}, y_{k,i}, \boldsymbol{\theta})\right)}{\int \exp\left(\sum_{k=1}^K \frac{-n_k}{n_0 \lambda_0} \frac{1}{n_k} \sum_{i=1}^{n_k} \ell_k(x_{k,i}, y_{k,i}, \mathbf{v})\right) dQ_1(\mathbf{v})} \quad (72)$$

$$= \frac{\exp\left(\frac{-1}{n_0 \lambda_0} \sum_{k=1}^K \sum_{i=1}^{n_k} \ell_k(x_{k,i}, y_{k,i}, \boldsymbol{\theta})\right)}{\int \exp\left(\frac{-1}{n_0 \lambda_0} \sum_{k=1}^K \sum_{i=1}^{n_k} \ell_k(x_{k,i}, y_{k,i}, \mathbf{v})\right) dQ_1(\mathbf{v})}. \quad (73)$$

where the equality in (72) follows from (12).

Under the assumption on the loss functions, *i.e.*, $\ell_1 = \ell_2 = \dots = \ell_K = \ell$ and the fact that \mathbf{z}_0 in (8) is the aggregation of $\mathbf{z}_1, \dots, \mathbf{z}_K$, and $n_0 = \sum_{k=1}^K n_k$, the equality in (73) yields for all $\boldsymbol{\theta} \in \text{supp } Q_1$,

$$\frac{dP_{\boldsymbol{\theta}_K | \mathbf{Z}_K = \mathbf{z}_K}^{(Q_K, \lambda_K)}(\boldsymbol{\theta})}{dQ_1} = \frac{\exp\left(\frac{-1}{n_0 \lambda_0} \sum_{i=1}^{n_0} \ell(x_{0,i}, y_{0,i}, \boldsymbol{\theta})\right)}{\int \exp\left(\frac{-1}{n_0 \lambda_0} \sum_{i=1}^{n_0} \ell(x_{0,i}, y_{0,i}, \mathbf{v})\right) dQ_1(\mathbf{v})} \quad (74)$$

$$= \frac{dP_{\boldsymbol{\theta}_K | \mathbf{Z}_0 = \mathbf{z}_0}^{(Q_1, \lambda_0)}(\boldsymbol{\theta})}{dQ_1}, \quad (75)$$

where the equality in (75) follows from (27).

Finally, from [20, Theorem 5] and (75), for all $\theta \in \text{supp } Q_1$,

$$1 = \frac{dP_{\Theta_K|Z_K=z_K}^{(Q_K, \lambda_K)}(\theta)}{dQ_1}(\theta) \frac{dQ_1}{dP_{\Theta_K|Z_0=z_0}^{(Q_1, \lambda_0)}}(\theta) \quad (76)$$

$$= \frac{dP_{\Theta_K|Z_K=z_K}^{(Q_K, \lambda_K)}(\theta)}{dP_{\Theta_K|Z_0=z_0}^{(Q_1, \lambda_0)}}(\theta), \quad (77)$$

where the equality in (77) follows from [20, Theorem 4] and completes the proof.

D Complementary Results

The following lemma provides an information-theoretic expression for C_k in (33) in terms of relative entropies.

Lemma 8. *The value C_k in (33) satisfies*

$$\begin{aligned} C_k &= D\left(P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)} \parallel Q_k\right) - D\left(P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)} \parallel Q_{k+1}\right) - D\left(P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)} \parallel Q_1\right) \\ &= D(Q_1 \parallel Q_k) - D(Q_1 \parallel Q_{k+1}) + D\left(Q_1 \parallel P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}\right), \end{aligned} \quad (79)$$

where the measure $P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}$, is an (L_k, Q_1, λ_k) -Gibbs probability measure; and the measure Q_k is defined in (21b).

Proof: From Lemma 7, for all $\theta \in \text{supp } Q_1$, it follows that

$$\frac{dP_{\Theta_k|Z_k=z_k}^{(Q_k, \lambda_k)}(\theta)}{dQ_k}(\theta) = \frac{dP_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}(\theta)}{dQ_1}(\theta) \exp(C_k), \quad (80)$$

with

$$C_k = \log \left(\frac{\int \exp\left(\sum_{i=1}^{k-1} \frac{-1}{\lambda_i} L_i(z_i, \tilde{\nu})\right) dQ_1(\tilde{\nu}) \left(\exp\left(K_{k, Q_1, z_k}\left(\frac{-1}{\lambda_k}\right)\right)\right)}{\int \exp\left(\sum_{j=1}^k \frac{-1}{\lambda_j} L_j(z_j, \nu)\right) dQ_1(\nu)} \right). \quad (81)$$

By taking a logarithm and integrating with respect to $P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}$ in (80), it follows that

$$\begin{aligned} & \int \log \left(\frac{dP_{\Theta_k|Z_k=z_k}^{(Q_k, \lambda_k)}(\theta)}{dQ_k}(\theta) \right) dP_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}(\theta) \\ &= \int \log \left(\frac{dP_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}(\theta)}{dQ_1}(\theta) \right) dP_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}(\theta) + C_k \end{aligned} \quad (82)$$

$$= D\left(P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)} \parallel Q_1\right) + C_k. \quad (83)$$

On the other hand, from Lemma 1 it follows that

$$\begin{aligned} & \int \log \left(\frac{dP_{\Theta_k | Z_k = z_k}^{(Q_k, \lambda_k)}}{dQ_k}(\theta) \right) dP_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)}(\theta) \\ &= D \left(P_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)} \parallel Q_k \right) - D \left(P_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_k, \lambda_k)} \right). \end{aligned} \quad (84)$$

Combining (83) and (84) yields

$$C_k = D \left(P_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)} \parallel Q_k \right) - D \left(P_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)} \parallel Q_{k+1} \right) - D \left(P_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)} \parallel Q_1 \right), \quad (85)$$

which completes the proof of (78).

Following the same method, if the integral in (80) is taken with respect to Q_1 it then follows that,

$$\begin{aligned} & \int \log \left(\frac{dP_{\Theta_k | Z_k = z_k}^{(Q_k, \lambda_k)}}{dQ_k}(\theta) \right) dQ_1(\theta) \\ &= \int \log \left(\frac{dP_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)}}{dQ_1}(\theta) \right) dQ_1(\theta) + C_k \end{aligned} \quad (86)$$

$$= \int \log \left(\left(\frac{dQ_1}{dP_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)}}(\theta) \right)^{-1} \right) dQ_1(\theta) + C_k \quad (87)$$

$$= - \int \log \left(\frac{dQ_1}{dP_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)}}(\theta) \right) dQ_1(\theta) + C_k \quad (88)$$

$$= -D \left(Q_1 \parallel P_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)} \right) + C_k, \quad (89)$$

where (87) follows from [6, Lemma 3], [20, Theorem 2] and [20, Theorem 5].

On the other hand, from Lemma 1 it follows that

$$\begin{aligned} & \int \log \left(\frac{dP_{\Theta_k | Z_k = z_k}^{(Q_k, \lambda_k)}}{dQ_k}(\theta) \right) dQ_1(\theta) \\ &= D(Q_1 \parallel Q_k) - D \left(Q_1 \parallel P_{\Theta_k | Z_k = z_k}^{(Q_k, \lambda_k)} \right). \end{aligned} \quad (90)$$

Combining (89) and (90) yields

$$C_k = D(Q_1 \parallel Q_k) - D(Q_1 \parallel Q_{k+1}) + D \left(Q_1 \parallel P_{\Theta_k | Z_k = z_k}^{(Q_1, \lambda_k)} \right), \quad (91)$$

which completes the proof of (79) and concludes the proof. \blacksquare

The explicit characterization of this constant makes clear that the nested choice of reference measures is not innocuous: it determines how the information contained in previous clients is incorporated through a client's reference measure

Q_k in (21b). This observation motivates comparing choices of reference measures in (21b) Q_1 or Q_k .

Lemma 8 implies a characterization of the Jeffreys divergence [22] between the (L_k, Q_1, λ_k) -Gibbs probability measure $P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}$ in (21) and Q_k in (21b), as shown by the following corollary.

Corollary 9. *The (L_k, Q_1, λ_k) -Gibbs probability measure $P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}$ in (21), satisfies*

$$\begin{aligned} & D\left(P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)} \parallel Q_1\right) + D\left(Q_1 \parallel P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)}\right) \\ = & D\left(P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)} \parallel Q_k\right) - D(Q_1 \parallel Q_k) \\ & + D(Q_1 \parallel Q_{k+1}) - D\left(P_{\Theta_k|Z_k=z_k}^{(Q_1, \lambda_k)} \parallel Q_{k+1}\right), \end{aligned} \quad (92)$$

where the measure Q_k is defined in (21b).



**RESEARCH CENTRE
SOPHIA ANTIPOLIS – MÉDITERRANÉE**

2004 route des Lucioles - BP 93
06902 Sophia Antipolis Cedex

Publisher
Inria
Domaine de Voluceau - Rocquencourt
BP 105 - 78153 Le Chesnay Cedex
inria.fr

ISSN 0249-6399