



HAL
open science

Statistical Federated Learning and Generalization

Yaiza Bermudez, Samir M. Perlaza, Iñaki Esnaola, H. Vincent Poor

► **To cite this version:**

Yaiza Bermudez, Samir M. Perlaza, Iñaki Esnaola, H. Vincent Poor. Statistical Federated Learning and Generalization. RR-9599, Inria. 2025. <hal-05355756>

HAL Id: hal-05355756

<https://inria.hal.science/hal-05355756v1>

Submitted on 13 Nov 2025

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



Statistical Federated Learning and Generalization

Yaiza Bermudez, Samir M. Perlaza, Iñaki Esnaola and
H. Vincent Poor

**RESEARCH
REPORT**

N° 9599

October 2025

Project-Team NEO

ISRN INRIA/RR--9599--FR+ENG

ISSN 0249-6399



Statistical Federated Learning and Generalization

Yaiza Bermudez, Samir M. Perlaza, Iñaki Esnaola and
H. Vincent Poor

Project-Team NEO

Research Report n° 9599 — October 2025 — 67 pages

Abstract: In this report, the generalization error of federated learning (FL) systems is characterized through a novel statistical framework. Central to this framework is the concept of a meta-federated learning algorithm, defined as a probability measure over a client’s local models conditioned on the datasets of all participating clients. By means of this abstraction, several fundamental properties of FL systems are stated and closed-form expressions for the generalization error are derived. More specifically, the method of gaps, originally introduced for non-federated settings, is extended to FL, and closed-form expressions for the generalization error are obtained in terms of classical information measures, including relative entropy, mutual information, and lautum information. A central role in these new expressions is played by some specific Gibbs probability measures (Gibbs algorithms). More importantly, it is revealed that the challenge of evaluating the generalization error in FL is reduced to two distinct tasks: (a) measuring the dependence of client model choices on the datasets of all clients; and (b) distinguishing the meta-federated learning algorithm from a Gibbs algorithm trained solely on local data. These findings establish new links between generalization in FL, mismatched hypothesis testing, Shannon’s information measures, and Pythagorean identities for the generalization error.

Key-words: Federated Learning, Generalization, Gibbs Measures, and Method of Gaps

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. H. Vincent Poor, Samir M. Perlaza, and Iñaki Esnaola 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 Fédéré et Généralisation

Résumé : Dans ce rapport, l'erreur de généralisation des systèmes d'apprentissage fédéré (FL) est caractérisée au moyen d'un nouveau cadre statistique. Au cœur de ce cadre se trouve la notion d'algorithme d'apprentissage méta-fédéré, défini comme une mesure de probabilité sur les modèles locaux d'un client, conditionnée par les ensembles de données de l'ensemble des clients participants. Grâce à cette abstraction, plusieurs propriétés fondamentales des systèmes de FL sont énoncées et des expressions explicites (en forme fermée) de l'erreur de généralisation sont dérivées. Plus précisément, la méthode des écarts (« method of gaps »), initialement introduite pour des contextes non fédérés, est étendue au FL, et des expressions en forme fermée de l'erreur de généralisation sont obtenues en termes de mesures d'information classiques, notamment l'entropie relative (divergence de Kullback-Leibler), l'information mutuelle et l'« information de lautum ». Un rôle central dans ces nouvelles expressions est joué par certaines mesures de probabilité de Gibbs (algorithmes de Gibbs). Plus important encore, il apparaît que la difficulté d'évaluer l'erreur de généralisation en FL se ramène à deux tâches distinctes : (a) mesurer la dépendance des choix de modèles des clients à l'égard des ensembles de données de tous les clients; et (b) distinguer l'algorithme d'apprentissage méta-fédéré d'un algorithme de Gibbs entraîné uniquement sur des données locales. Ces résultats établissent de nouveaux liens entre la généralisation en FL, les tests d'hypothèses « mismatched » (non adaptés), les mesures d'information de Shannon et des identités pythagoriciennes pour l'erreur de généralisation.

Mots-clés : Apprentissage Fédéré, Généralisation, Mesures de Gibbs et Méthode des Ecarts

Contents

1	Introduction	5
1.1	Related Work	5
1.2	Contributions	5
2	System Model	6
2.1	Communication System	8
2.1.1	Preamble	8
2.1.2	Communication Protocol	8
2.2	Stochastic Modelling	9
2.2.1	Model of the Data Source	13
2.2.2	Model of the Server	14
2.2.3	Model of the Clients	16
3	Generalization Error	18
3.1	Meta-Federated Learning Algorithm	18
3.2	Expected Empirical Risks	19
3.3	Definition of Generalization Error	21
4	Gibbs Measures	23
4.1	Data-Dependent Gibbs Measures	23
4.2	Model-Dependent Gibbs Measures	25
5	Method of Gaps	27
5.1	Expected Empirical Risk Gaps	27
5.2	Step One	29
5.3	Step Two	29
6	Data-dependent Closed-Form Expressions	31
6.1	The First Closed-Form Expression	31
6.2	The Second Closed-Form Expression	32
6.3	A Pythagorean Identity of Generalization Error	33
7	Model-dependent Closed-Form Expressions	34
7.1	The First Closed-Form Expression	34
7.2	The Second Closed-Form Expression	36
7.3	A Pythagorean Identity of Generalization Error	36
	References	37
A	Proof of Lemma 1	44
B	Proof of Lemma 2	45
C	Proof of Lemma 3	46
D	Proof of Lemma 4	47

E	Proof of Lemma 5	48
F	Proof of Lemma 12	48
G	Proof of Lemma 13	50
H	Proof of Lemma 14	51
I	Proof of Lemma 15	52
J	Proof of Theorem 18	53
K	Proof of Theorem 19	54
L	Proof of Theorem 20	57
M	Proof of Theorem 21	58
N	Proof of Theorem 22	60
O	Proof of Theorem 23	62
P	Proof of Theorem 24	64
Q	Proof of Theorem 25	66

1 Introduction

Federated learning (FL), first introduced in [1], harnesses training across a population of clients by repeatedly exchanging information with a coordinating server [2, 3]. This setting departs from classical centralized learning in two crucial ways: (i) data are statistically heterogeneous across clients [4–7] and often severely unbalanced in terms of integrity and dataset sizes [1, 2]; and (ii) the set of participating clients varies over time due to system constraints, communication costs, and availability [8–10].

From this perspective, the analysis of generalization capabilities of FL has been out of the reach of classical methods based on VC-dimension or Rademacher complexity [11–15], despite important contributions, such as [4, 16, 17]. This contrasts with non-federated learning for which a rich variety of bounds [18–27] and exact expressions are known for broad classes of statistical learning algorithms. Most of the existing exact expressions for the generalization error were separately introduced by different authors using specific methods of proofs for particular learning algorithms [26, 28, 29]. Interestingly, the method of gaps, introduced in [30], establishes a unified framework from which all existing exact expressions can be obtained. More importantly, such a method establishes strong connections among generalization and several areas in statistics, e.g., mismatched hypothesis testing [31], classical information measures [32–34], Pythagorean identities of variations of expectations [35], and I-projections [36].

1.1 Related Work

Optimization and statistical methods have led to insightful lower and upper bounds on the generalization error of FL systems [9, 17, 37–50]. Some contributions focus on a *participation* (client-sampling) generalization gap, i.e., the discrepancy between the generalization performance obtained when training involves a random subset of clients and the performance exclusively involving local training datasets. Along the participation generalization gap, also an *out-of-sample* generalization gap has been studied in [51, 52], where out-of-sample generalization must be understood in the sense of [53]. The interest in these gaps lies in the fact that their sum equals the generalization error of the FL system under study. These gaps are typically studied via optimization-centric convergence analyses for specific algorithms, including FedAvg [1], FedProx [5], SCAFFOLD [6], FedOpt [54], among others, under heterogeneous data and partial participation [55, 56]. Despite the numerous lower and upper bounds on the generalization error of federated learning algorithms, exact expressions are uncharted territory.

1.2 Contributions

This report develops a statistical model that explicitly characterizes, in statistical terms, each of the components of an FL system. For a fixed set of participating clients, this model reproduces as special cases, most standard FL scenarios,

e.g., FedAvg [1], FedProx [5], among others. Central to this formulation is a new notion, referred to as *meta-federated learning algorithm*, which consists of a probability measure over each client’s local models conditioned on the joint training datasets of all clients. This conditional probability measure arises naturally from the client-server communication process and serves as a stand-alone mathematical abstraction of such a client within the federation. Each client selects its models by sampling such a probability measure. Leveraging this flexible and general abstraction, important conclusions are drawn for all FL systems. First, it is shown that, the meta-federated learning algorithm of each client is a convex combination of all instances of the local algorithm, which is represented by a probability measure on the models conditioned on the local training dataset and all messages exchanged with the server. Hence, the set of all meta-federated learning algorithms in FL is a strict subset of all probability measures that can be defined over the set of models. Second, it is shown that relevant performance metrics for a given client, such as the expected test empirical risk, the expected training empirical risk, and the expected generalization error, can be expressed in terms of the corresponding meta-federated learning algorithm. This observation enables an extension of the method of gaps, introduced in [30] for the case of non-federated systems, to the federated setting. As a result, closed-form expressions for the generalization error of FL systems are obtained in terms of classical information measures (relative entropy, mutual information [32, 33], and lautum information [34]). In such expressions, a Gibbs probability measure on the models, which represents a Gibbs algorithm [57], plays a central role. In particular, the difficulty of computing a client’s generalization error in an FL system is shown to be decomposed into two separated challenges: (i) quantifying the dependence of the client’s model selection on the training datasets of all participating clients via the mutual and lautum information induced by the meta-federated learning algorithm; and (ii) reliably distinguishing the meta-federated learning algorithm from a Gibbs algorithm exclusively trained on the individual training dataset. These challenges unveil exciting connections among generalization in FL, mismatched hypothesis testing, Shannon’s information measures and Pythagorean identities for the generalization error.

2 System Model

Consider a federated learning system in which K clients collaboratively tune their local learning algorithms by communicating with a common server. 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

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

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

$$\mathbf{z}_0 \triangleq (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K) \in \mathcal{Z}_1^{n_1} \times \mathcal{Z}_2^{n_2} \times \dots \times \mathcal{Z}_K^{n_K} \quad (2)$$

$$= ((x_{0,1}, y_{0,1}), (x_{0,2}, y_{0,2}), \dots, (x_{0,n_0}, y_{0,n_0})). \quad (3)$$

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

$$n_0 \triangleq \sum_{k=1}^K n_k. \quad (4)$$

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

$$\ell_k : \mathcal{Z}_k \times \mathcal{M}_k \rightarrow [0, +\infty), \quad (5)$$

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 (1), is determined by the function

$$\mathbb{L}_k : \begin{cases} \mathcal{Z}_k^{n_k} \times \mathcal{M}_k \longrightarrow [0, +\infty) \\ (\mathbf{z}_k, \theta) \longmapsto \frac{1}{n_k} \sum_{i=1}^{n_k} \ell_k(x_{k,i}, y_{k,i}, \theta), \end{cases} \quad (6)$$

where the function ℓ_k is defined in (5).

The underlying assumption of federated learning is that each client sends to and receives messages from a server. Such messages, which circulate through a communication network, help clients to cooperatively search for models. A typical formulation of a federated learning system consists in distributively solving the following optimization problem via messages exchanges between the server and the clients,

$$\min_{\theta \in \mathcal{M}} \sum_{k=1}^K \frac{n_k}{n_0} \mathbb{L}_k(\mathbf{z}_k, \theta), \quad (7)$$

where $\mathcal{M} \triangleq \bigcap_{i=1}^K \mathcal{M}_k$ is the set of common models. The motivation of clients for adopting models that are solutions to the problem in (7) becomes evident in the case in which all agents use the same loss function, i.e., $\ell_1 = \ell_2 = \dots = \ell_K = \ell$. In such a case, the problem in (7) becomes

$$\min_{\theta \in \mathcal{M}} \frac{1}{n_0} \sum_{t=1}^{n_0} \ell(\theta, x_{0,t}, y_{0,t}), \quad (8)$$

where the pair $(x_{0,t}, y_{0,t})$ is a data point in the dataset \mathbf{z}_0 in (3). Note such a dataset \mathbf{z}_0 aggregates all local datasets and thus, the solution to (8) can be seen as a model that minimizes the empirical risk with respect to the aggregation of all local datasets.

In the following, the communication network is described. Using such a network, the stochastic description of the federated learning system under study is presented. Finally, the mathematical formulation of the problem studied in this work is thoroughly described.

2.1 Communication System

Server and clients engage on a synchronous communication over bidirectional, noiseless channels with information transmission rate limitations. This communication consists of m communication rounds. In each round, all clients send a message to the server, which subsequently sends a message to each client. At the end of these m communication rounds, all clients locally choose their corresponding models. Previous to this communication, a preamble is conducted to configure both local and global parameters for the clients and the server.

2.1.1 Preamble

The preamble precedes the federated learning process and serves to initialize the system. During this period, clients acquire reference probability measures, typically in the form of priors, over the set of models and/or the set of datasets. Such reference measures may be transmitted by the server or derived locally by each client using available data, e.g., training datasets.

2.1.2 Communication Protocol

During communication round t , with $t \in \{1, 2, \dots, m\}$ and $k \in \{1, 2, \dots, K\}$, the message received by client k from the server and the message sent by client k to the server are denoted by

$$u_{k,t} \in \mathcal{U}_k, \text{ and} \tag{9}$$

$$v_{k,t} \in \mathcal{V}_k, \tag{10}$$

respectively, with \mathcal{U}_k and \mathcal{V}_k some given finite sets. Communication round t starts with all clients sending their corresponding messages $v_{1,t}, v_{2,t}, \dots, v_{K,t}$ to the server. Once all messages have been received by the server, it sends the messages $u_{1,t}, u_{2,t}, \dots, u_{K,t}$ to the corresponding clients. Using this notation, Figure 1 depicts the exchange of messages during two communication rounds, namely $t - 1$ and t , in the case in which $K = 2$ and $t > 1$.

For all $t \in \{1, 2, \dots, m\}$, messages $v_{1,t}, v_{2,t}, \dots, v_{K,t}$ and $u_{1,t}, u_{2,t}, \dots, u_{K,t}$ are subject to information transmission rate constraints. More specifically, at each communication round, at most $\bar{R}_k \in (0, +\infty)$ bits are sent by client k to the server; and at most $\underline{R}_k \in (0, +\infty)$ bits are received by client k from the server. Hence, \bar{R}_k and \underline{R}_k are the maximum information transmission rates in bits per

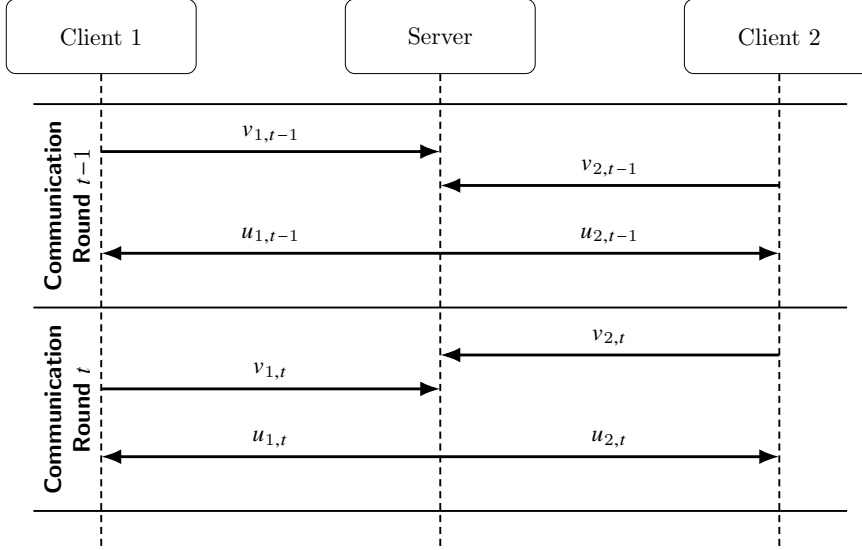


Figure 1: Two consecutive communication rounds, namely $t - 1$ and t , with $t \in \{2, 3, \dots, m\}$, in the case in which $K = 2$.

communication-round duration. These constraints translate into the following inequalities, for all $k \in \{1, 2, \dots, K\}$,

$$|\mathcal{U}_k| \leq 2^{\underline{R}_k} \text{ and} \quad (11)$$

$$|\mathcal{V}_k| \leq 2^{\bar{R}_k}. \quad (12)$$

This is due to the fact that for establishing a one-to-one relation between the elements of the set \mathcal{U}_k (respectively \mathcal{V}_k) and some binary sequences, it requires having $|\mathcal{U}_k|$ (respectively $|\mathcal{V}_k|$) sequences of at least $\log_2(|\mathcal{U}_k|)$ bits (respectively $\log_2(|\mathcal{V}_k|)$ bits). This implies that the rate \underline{R}_k (respectively \bar{R}_k) shall be at least $\log_2(|\mathcal{U}_k|)$ (respectively $\log_2(|\mathcal{V}_k|)$) bits per communication round from the server to client k (respectively from client k to the server). This justifies the inequalities (11) and (12). In the following, the uplink and downlink data rates, \bar{R}_k and \underline{R}_k respectively, are fixed parameters in this analysis.

2.2 Stochastic Modelling

For the ease of notation, all messages sent to the server by the clients during round t are denoted by

$$\mathbf{v}^{(t)} = (v_{1,t}, v_{2,t}, \dots, v_{K,t})^\top \in \mathcal{V}_1 \times \mathcal{V}_2 \times \dots \times \mathcal{V}_K; \quad (13)$$

and all messages sent by client k during all previous rounds, including round t , are denoted by

$$\mathbf{v}_k^{(t)} = (v_{k,1}, v_{k,2}, \dots, v_{k,t})^\top \in \mathcal{V}_k^t. \quad (14)$$

The messages sent by the server to the clients, during such a communication round t , are denoted by

$$\mathbf{u}^{(t)} = (u_{1,t}, u_{2,t}, \dots, u_{K,t})^\top \in \mathcal{U}_1 \times \mathcal{U}_2 \times \dots \times \mathcal{U}_K; \quad (15)$$

and all the messages sent by the server to client k during all previous rounds, including round t , are denoted by

$$\mathbf{u}_k^{(t)} = (u_{k,1}, u_{k,2}, \dots, u_{k,t})^\top \in \mathcal{U}_k^t. \quad (16)$$

Note that the tuples in (13) and (14) share the following structure,

$$\left(\mathbf{v}_1^{(t)} \ \mathbf{v}_2^{(t)} \ \dots \ \mathbf{v}_K^{(t)} \right)^\top = \underbrace{\begin{bmatrix} v_{1,1} & v_{1,2} & \dots & v_{1,t} \\ v_{2,1} & v_{2,2} & \dots & v_{2,t} \\ \vdots & \vdots & \ddots & \vdots \\ v_{K,1} & v_{K,2} & \dots & v_{K,t} \end{bmatrix}}_{\triangleq \mathbf{v}^{(t)} \in (\mathcal{V}_1 \times \mathcal{V}_2 \times \dots \times \mathcal{V}_K)^t} = \left(\mathbf{v}^{(1)} \ \mathbf{v}^{(2)} \ \dots \ \mathbf{v}^{(t)} \right). \quad (17)$$

Similarly, the tuples in (15) and (16) share the following structure,

$$\left(\mathbf{u}_1^{(t)} \ \mathbf{u}_2^{(t)} \ \dots \ \mathbf{u}_K^{(t)} \right)^\top = \underbrace{\begin{bmatrix} u_{1,1} & u_{1,2} & \dots & u_{1,t} \\ u_{2,1} & u_{2,2} & \dots & u_{2,t} \\ \vdots & \vdots & \ddots & \vdots \\ u_{K,1} & u_{K,2} & \dots & u_{K,t} \end{bmatrix}}_{\triangleq \mathbf{u}^{(t)} \in (\mathcal{U}_1 \times \mathcal{U}_2 \times \dots \times \mathcal{U}_K)^t} = \left(\mathbf{u}^{(1)} \ \mathbf{u}^{(2)} \ \dots \ \mathbf{u}^{(t)} \right). \quad (18)$$

For the ease of notation, let the sets \mathcal{M}_0 , \mathcal{U}_0 , \mathcal{V}_0 , and \mathcal{Z}_0 be respectively defined as follows:

$$\mathcal{M}_0 \triangleq \mathcal{M}_1 \times \mathcal{M}_2 \times \dots \times \mathcal{M}_K \subset \mathbb{R}^{d \times K}; \quad (19)$$

$$\mathcal{U}_0 \triangleq \mathcal{U}_1 \times \mathcal{U}_2 \times \dots \times \mathcal{U}_K; \quad (20)$$

$$\mathcal{V}_0 \triangleq \mathcal{V}_1 \times \mathcal{V}_2 \times \dots \times \mathcal{V}_K; \quad (21)$$

$$\mathcal{Z}_0 \triangleq \mathcal{Z}_1^{n_1} \times \mathcal{Z}_2^{n_2} \times \dots \times \mathcal{Z}_K^{n_K}. \quad (22)$$

After m communication rounds, the model chosen by client k is denoted by $\theta_k \in \mathcal{M}_k$. The tuple containing the models chosen by all clients is denoted by

$$\underline{\theta} = (\theta_1, \theta_2, \dots, \theta_K) \in \mathcal{M}_0. \quad (23)$$

In this analysis, at the end of m communication rounds, the dataset, \mathbf{z}_0 in (2); the messages sent and received throughout the communication process, $\underline{\mathbf{v}}^{(m)}$

in (17) and $\underline{\mathbf{u}}^{(m)}$ and (18); and the models chosen by all clients, $\underline{\boldsymbol{\theta}}$ in (23), are obtained by sampling the joint probability measure

$$P_{\underline{\boldsymbol{\theta}}, \underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, \mathbf{z}_0} \in \Delta(\mathcal{M}_0 \times \mathcal{U}_0^m \times \mathcal{V}_0^m \times \mathcal{Z}_0), \quad (24)$$

where the sets \mathcal{M}_0 , \mathcal{U}_0 , \mathcal{V}_0 , and \mathcal{Z}_0 are defined in (19), (20), (21), and (22), respectively. The joint probability measure in (24) can be factorized as

$$P_{\underline{\boldsymbol{\theta}}, \underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, \mathbf{z}_0} = P_{\underline{\boldsymbol{\theta}}, \underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)} | \mathbf{z}_0} P_{\mathbf{z}_0}, \quad (25)$$

with,

$$P_{\underline{\boldsymbol{\theta}}, \underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)} | \mathbf{z}_0} \in \Delta(\mathcal{M}_0 \times \mathcal{U}_0^m \times \mathcal{V}_0^m | \mathcal{Z}_0); \text{ and} \quad (26)$$

$$P_{\mathbf{z}_0} \in \Delta(\mathcal{Z}_0). \quad (27)$$

The conditional probability measure $P_{\underline{\boldsymbol{\theta}}, \underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)} | \mathbf{z}_0}$ in (26) is a statistical description of the federated learning system. Alternatively, the probability measure $P_{\mathbf{z}_0}$ in (27) is a statistical description of the data source.

The analysis continues by highlighting some of the constraints derived from the federated learning system implementation.

Assumption 1. *The model chosen by client k , with $k \in \{1, 2, \dots, K\}$, is obtained by sampling a probability measure exclusively conditioned on*

- The dataset \mathbf{z}_k in (1);
- The m messages sent to the server, $\mathbf{v}_k^{(m)}$ in (14); and
- The m messages received from the server, $\mathbf{u}_k^{(m)}$ in (16).

The conditional probability measure described in Assumption 1 is denoted by

$$P_{\boldsymbol{\theta}_k | \mathbf{u}_k^{(m)}, \mathbf{v}_k^{(m)}, \mathbf{z}_k} \in \Delta(\mathcal{M}_k | \mathcal{U}_k^m \times \mathcal{V}_k^m \times \mathcal{Z}_k^{n_k}), \quad (28)$$

and can be interpreted as the *machine learning algorithm* used by client k , whose inputs are the dataset \mathbf{z}_k in (1); the m messages sent to the server, $\mathbf{v}_k^{(m)}$ in (14); and the m messages received from the server, $\mathbf{u}_k^{(m)}$ in (16).

Assumption 2. *The message sent by client k to the server at communication round t , with $k \in \{1, 2, \dots, K\}$ and $t \in \{1, 2, \dots, m\}$, is obtained by sampling a probability measure exclusively conditioned on*

- The dataset \mathbf{z}_k in (1);
- The $t - 1$ messages previously sent to the server, $\mathbf{v}_k^{(t-1)}$ in (14); and
- The $t - 1$ messages previously received from the server, $\mathbf{u}_k^{(t-1)}$ in (16).

The conditional probability measure described in Assumption 2, at communication round t , with $t \in \{1, 2, \dots, m\}$, is denoted by

$$P_{V_{k,t}|\underline{U}_k^{(t-1)}, \underline{V}_k^{(t-1)}, \underline{Z}_k} \in \Delta \left(\mathcal{V}_k | \mathcal{U}_k^{t-1} \times \mathcal{V}_k^{t-1} \times \mathcal{Z}_k^{n_k} \right), \quad (29)$$

and can be interpreted as a mathematical abstraction of the part of client k that determines the message to be sent to the server during communication round t .

Assumption 3. *The messages sent by the server at communication round t to all clients, $\mathbf{u}^{(t)}$ in (15), with $t \in \{1, 2, \dots, m\}$, are obtained by sampling a probability measure exclusively conditioned on*

- *The t messages received from each client, $\underline{\mathbf{v}}^{(t)}$ in (17); and*
- *The $t - 1$ messages previously sent to the clients, $\underline{\mathbf{u}}^{(t-1)}$ in (18).*

The conditional probability measure described in Assumption 3, at communication round t , with $t \in \{1, 2, \dots, m\}$, is denoted by

$$P_{\underline{U}^{(t)}|\underline{U}^{(t-1)}, \underline{V}^{(t)}} \in \Delta \left(\mathcal{U}_0 | \mathcal{U}_0^{t-1} \times \mathcal{V}_0^t \right), \quad (30)$$

and can be seen as a mathematical abstraction of the server during communication round t . This follows by observing that the messages sent to the clients during communication round t , $\mathbf{u}^{(t)}$ in (15), are chosen exclusively based on all messages previously sent to the clients $\mathbf{u}^{(t-1)}$; and on all messages previously received from the clients, including those of the current communication round, $\mathbf{v}^{(t)}$ in (13), which is all the data available at the server.

The following lemma puts the above assumptions at play to factorize the conditional probability measure $P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|\underline{Z}_0}$ in (26).

Lemma 1. *Under Assumptions 1, 2, and 3, the joint probability measure $P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|\underline{Z}_0}$ in (26) satisfies,*

$$P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|\underline{Z}_0} = \left(\prod_{k=1}^K P_{\Theta_k | U_k^{(m)}, V_k^{(m)}, Z_k} \right) \prod_{t=1}^m \left(P_{\underline{U}^{(t)}|\underline{U}^{(t-1)}, \underline{V}^{(t)}} \prod_{k=1}^K P_{V_{k,t} | U_k^{(t-1)}, V_k^{(t-1)}, Z_k} \right), \quad (31)$$

where the conditional probability measure $P_{\Theta_k | U_k^{(m)}, V_k^{(m)}, Z_k}$ is defined in (28); the conditional probability measure $P_{V_{k,t} | U_k^{(t-1)}, V_k^{(t-1)}, Z_k}$ is defined in (29); and the conditional probability measure $P_{\underline{U}^{(t)}|\underline{U}^{(t-1)}, \underline{V}^{(t)}}$ is defined in (30).

Proof: The proof is presented in Appendix A. ■

Essentially, given the m exchanged messages, $\underline{\mathbf{v}}^{(m)}$ and $\underline{\mathbf{u}}^{(m)}$ in (13) and (15); and the local training dataset \mathbf{z}_k , client k chooses its models by sampling the probability measure $P_{\Theta_k | U_k^{(m)} = \mathbf{u}_k^{(m)}, V_k^{(m)} = \mathbf{v}_k^{(m)}, Z_k = \mathbf{z}_k} \in \Delta(\mathcal{M}_k)$. The message sent by

client k to the server at communication round t , $v_{k,t}$ in (10), is obtained by sampling the probability measure $P_{V_{k,t} | \mathbf{U}_k^{(t-1)} = \mathbf{u}_k^{(t-1)}, \mathbf{V}_k^{(t-1)} = \mathbf{v}_k^{(t-1)}, \mathbf{Z}_k = \mathbf{z}_k} \in \Delta(\mathcal{V}_k)$. Alternatively, the messages sent by the server at communication round t to the clients can be obtained by sampling the probability measure $P_{\mathbf{U}^{(t)} | \mathbf{U}^{(t-1)} = \mathbf{u}^{(t-1)}, \mathbf{V}^{(t)} = \mathbf{v}^{(t)}} \in \Delta(\mathcal{U}_0)$. Using these measures, it is possible to establish mathematical representations for all objects intervening in the learning process, namely, the data source, the server, and the clients. See for instance, Figure 2, which describes a *one-shot* federated learning, as is often referred to the case in which $m = 1$.

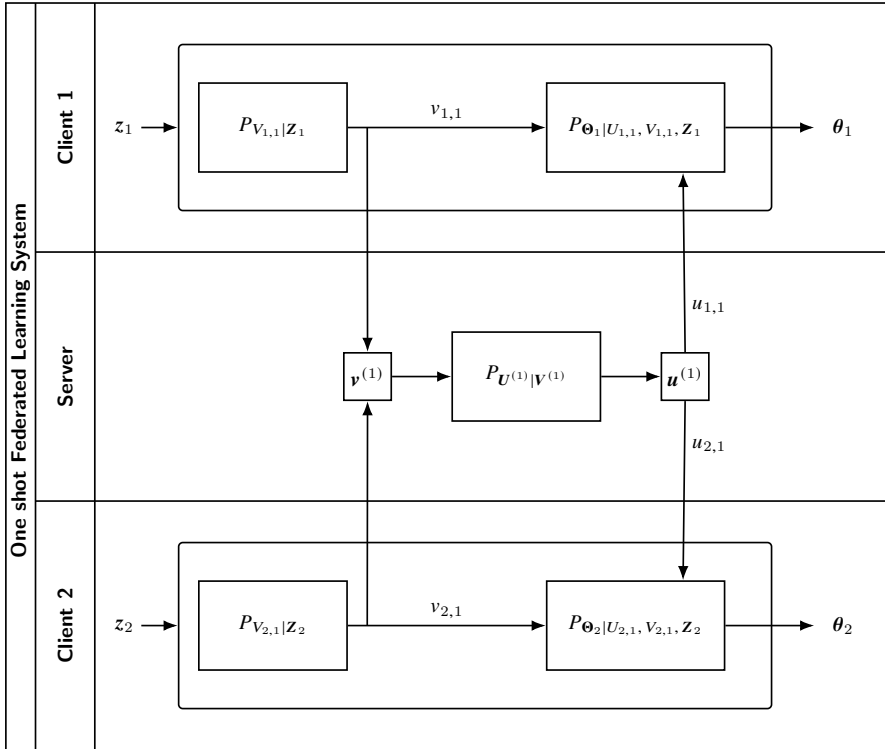


Figure 2: System model with $K = 2$ and $m = 1$ (one-shot federated learning).

2.2.1 Model of the Data Source

The data source is modeled by the probability measure $P_{\mathbf{Z}_0}$ in (27). That is, the datasets $\mathbf{z}_0 = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K)$ in (2) are obtained by sampling such a measure $P_{\mathbf{Z}_0}$. The marginal probability measures of $P_{\mathbf{Z}_0}$ in $\Delta(\mathcal{Z}_k^{n_k})$, with $k \in \{1, 2, \dots, K\}$, are denoted by $P_{\mathbf{Z}_k}$ and are defined as follows. For all $k \in \{1, 2, \dots, K\}$ and for

all measurable sets $\mathcal{A} \subset \mathcal{Z}_k^{n_k}$,

$$P_{\mathbf{Z}_k}(\mathcal{A}) = P_{\mathbf{Z}_0} \left(\prod_{i=1}^{k-1} \mathcal{Z}_i^{n_i} \times \mathcal{A} \times \prod_{i=k+1}^K \mathcal{Z}_i^{n_i} \right). \quad (32)$$

A common assumption on data sources is that they generate independent and identically distributed datasets. The following definition formalizes these assumptions.

Definition 1 (Independent and Identically Distributed Data Sets). *A data source represented by the probability measure $P_{\mathbf{Z}_0}$ in (27) is said to generate independent datasets if for all measurable sets $\mathcal{A} = \mathcal{A}_1 \times \mathcal{A}_2 \times \dots \times \mathcal{A}_K$, with $\mathcal{A}_k \subset \mathcal{Z}_k^{n_k}$,*

$$P_{\mathbf{Z}_0}(\mathcal{A}) = \prod_{k=1}^K P_{\mathbf{Z}_k}(\mathcal{A}_k), \quad (33)$$

where the probability measures $P_{\mathbf{Z}_1}, \dots, P_{\mathbf{Z}_K}$ are defined in (32). Moreover, such datasets are said to be identically distributed if

$$P_{\mathbf{Z}_1} = P_{\mathbf{Z}_2} = \dots = P_{\mathbf{Z}_K}. \quad (34)$$

Another common assumption consists in considering that datasets are formed by independent datapoints. The following definition formalizes such an assumption.

Definition 2 (Independent and Identically Distributed Data Points). *A data source represented by the probability measure $P_{\mathbf{Z}_0}$ in (27) is said to generate datasets formed by independent and identically distributed data points, if for all $k \in \{1, 2, \dots, K\}$, the marginal measure $P_{\mathbf{Z}_k}$ in (32) is a product measure formed by a probability measure $P_{\mathbf{Z}_k} \in \Delta(\mathcal{Z}_k)$. That is, for all measurable sets $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_{n_k}$, with $\mathcal{A}_j \subset \mathcal{Z}_k$, and $j \in \{1, 2, \dots, n_k\}$, it holds that*

$$P_{\mathbf{Z}_k}(\mathcal{A}) = \prod_{j=1}^{n_k} P_{\mathbf{Z}_k}(\mathcal{A}_j). \quad (35)$$

2.2.2 Model of the Server

At each communication round, the server is modeled by a conditional probability measure, from which the messages sent to the clients are sampled. During communication round t , with $t \in \{1, 2, \dots, m\}$, the server is modeled by the conditional probability measure $P_{\mathbf{U}^{(t)} | \mathbf{V}^{(t)}, \mathbf{U}^{(t-1)}}$ in (30). The following definitions introduce the most common assumptions on the server.

Definition 3 (Deterministic servers). *A server is said to be deterministic if there exists a family of functions*

$$g_t : \mathcal{V}_0^t \times \mathcal{U}_0^{t-1} \rightarrow \mathcal{U}_0, \quad \text{with } t \in \{1, 2, \dots, m\}, \quad (36)$$

such that for all $\underline{\mathbf{u}}^{(t-1)} \in \mathcal{U}_0^{t-1}$, for all $\underline{\mathbf{v}}^{(t)} \in \mathcal{V}_0^t$, and for all measurable sets $\mathcal{A} \subseteq \mathcal{U}_0$, the conditional probability measure $P_{\mathbf{U}^{(t)}|\underline{\mathbf{V}}^{(t)},\underline{\mathbf{U}}^{(t-1)}}$ in (30) satisfies

$$P_{\mathbf{U}^{(t)}|\underline{\mathbf{V}}^{(t)}=\underline{\mathbf{v}}^{(t)},\underline{\mathbf{U}}^{(t-1)}=\underline{\mathbf{u}}^{(t-1)}}(\mathcal{A}) = \begin{cases} 1 & \text{if } g_t(\underline{\mathbf{v}}^{(t)}, \underline{\mathbf{u}}^{(t-1)}) \in \mathcal{A}, \\ 0 & \text{otherwise.} \end{cases} \quad (37)$$

Essentially, the messages sent to the clients by a deterministic server are functions of all messages previously received from the clients, including those received during the current communication round, as well as, of all messages previously sent to the clients. That is, the following holds with probability one (w.r.t. $P_{\mathbf{U}^{(t)}|\underline{\mathbf{V}}^{(t)}=\underline{\mathbf{v}}^{(t)},\underline{\mathbf{U}}^{(t-1)}=\underline{\mathbf{u}}^{(t-1)}}$ in (37), with $(\underline{\mathbf{v}}^{(t)}, \underline{\mathbf{u}}^{(t-1)}) \in \mathcal{V}_0^t \times \mathcal{U}_0^{t-1}$):

$$\mathbf{u}^{(1)} = g_1(\mathbf{v}^{(1)}) \quad \text{and} \quad (38)$$

$$\mathbf{u}^{(t)} = g_t(\mathbf{v}^{(t)}, \mathbf{u}^{(t-1)}), \quad \text{with } t \in \{2, 3, \dots, m\}. \quad (39)$$

Definition 4 (Memoryless server). *A server is said to be memoryless if for all $t \in \{2, 3, \dots, m\}$, the conditional probability measure $P_{\mathbf{U}^{(t)}|\underline{\mathbf{V}}^{(t)},\underline{\mathbf{U}}^{(t-1)}}$ in (30) satisfies*

$$P_{\mathbf{U}^{(t)}|\underline{\mathbf{V}}^{(t)},\underline{\mathbf{U}}^{(t-1)}} = P_{\mathbf{U}^{(t)}|\mathbf{V}^{(t)}}, \quad (40)$$

for some given conditional probability measures

$$P_{\mathbf{U}^{(1)}|\mathbf{V}^{(1)}}, \quad P_{\mathbf{U}^{(2)}|\mathbf{V}^{(2)}}, \quad \dots, \quad P_{\mathbf{U}^{(m)}|\mathbf{V}^{(m)}} \text{ in } \Delta(\mathcal{U}_0|\mathcal{V}_0).$$

At each communication round t , with $t \in \{1, 2, 3, \dots, m\}$, the messages sent by a memoryless server depend exclusively on the messages received by the server at the beginning of the current communication round. Servers that are not memoryless, are said to be servers *with memory*.

Definition 5 (Time-Invariant server). *A server is said to be time-invariant if for all $t \in \{2, 3, \dots, m\}$, the conditional probability measure $P_{\mathbf{U}^{(t)}|\underline{\mathbf{V}}^{(t)},\underline{\mathbf{U}}^{(t-1)}}$ in (30) satisfies*

$$P_{\mathbf{U}^{(1)}|\mathbf{V}^{(1)}} = P_{\mathbf{U}^{(t)}|\underline{\mathbf{V}}^{(t)},\underline{\mathbf{U}}^{(t-1)}} = P_{\mathbf{U}|\mathbf{V}}, \quad (41)$$

for some fixed conditional probability measure $P_{\mathbf{U}|\mathbf{V}} \in \Delta(\mathcal{U}_0|\mathcal{V}_0)$.

Servers that are time-invariant are also memoryless. Servers that are not time-invariant are referred to as *time-varying* servers.

Definition 6 (Uniform server). *A server is said to be uniform if, at each communication round, it sends identical messages to all clients. That is, for all $t \in \{1, 2, \dots, m\}$,*

$$u_{1,t} = u_{2,t} = \dots = u_{K,t}, \quad (42)$$

where $u_{1,t}, u_{2,t}, \dots, u_{K,t}$ are the individual messages sent to the clients during round t , as defined in (9).

Uniform servers are also referred to as *homogeneous, broadcast, or global* servers. Servers that are not uniform are often called *heterogeneous or personalized* servers.

2.2.3 Model of the Clients

A client is modeled by two objects, namely the *machine learning algorithm* and the *message generator* or *messenger*. At Client k , with $k \in \{1, 2, \dots, K\}$, the former is modeled by the conditional probability measure $P_{\Theta_k|U_k^{(m)}, V_k^{(m)}, Z_k}$ in (28). The latter is modeled by the measures $P_{V_{k,t}|U_k^{(t-1)}, V_k^{(t-1)}, Z_k}$ in (29), with $t \in \{1, 2, \dots, m\}$. The following definitions describe the main assumptions on a client expressed in terms of such a client's machine learning algorithm and messenger.

Definition 7 (Deterministic Client). *Client k , with $k \in \{1, 2, \dots, K\}$, is said to be deterministic if there exists a family of functions*

$$g_{k,t} : \mathcal{U}_k^{t-1} \times \mathcal{V}_k^{t-1} \times \mathcal{Z}_k^{n_k} \rightarrow \mathcal{V}_k, \text{ with } t \in \{1, 2, \dots, m\}, \quad (43)$$

such that for all $\mathbf{u}_k^{(t-1)} \in \mathcal{U}_k^{(t-1)}$, for all $\mathbf{v}_k^{(t-1)} \in \mathcal{V}_k^{(t-1)}$, for all $\mathbf{z}_k \in \mathcal{Z}_k$, and for all measurable sets $\mathcal{B} \in \mathcal{V}_k$, the conditional probability measure $P_{V_{k,t}|U_k^{(t-1)}, V_k^{(t-1)}, Z_k}$ in (29) satisfies

$$P_{V_{k,t}|U_k^{(t-1)}=\mathbf{u}_k^{(t-1)}, V_k^{(t-1)}=\mathbf{v}_k^{(t-1)}, Z_k=\mathbf{z}_k}(\mathcal{B}) = \begin{cases} 1 & \text{if } g_{k,t}(\mathbf{u}_k^{(t-1)}, \mathbf{v}_k^{(t-1)}, \mathbf{z}_k) \in \mathcal{B} \\ 0 & \text{otherwise.} \end{cases} \quad (44)$$

Equation (44) implies, that the messenger is deterministic, *i.e.* operates without randomness in the generation of the transmitted messages. A client is therefore deterministic when its messenger is deterministic.

Definition 8 (Conforming Client). *Client k , with $k \in \{1, 2, \dots, K\}$, is said to be conforming if*

1. *the conditional probability measure $P_{\Theta_k|U_k^{(m)}, V_k^{(m)}, Z_k}$ in (28) satisfies*

$$P_{\Theta_k|U_k^{(m)}, V_k^{(m)}, Z_k} = P_{\Theta_k|U_k^{(m)}, Z_k}, \quad (45)$$

for some conditional probability measure $P_{\Theta_k|U_k^{(m)}, Z_k} \in \Delta(\mathcal{M}_k|U_k^m \times \mathcal{Z}_k^{n_k})$;
and

2. *for all $t \in \{1, 2, \dots, m\}$, the conditional probability measure $P_{V_{k,t}|U_k^{(t-1)}, V_k^{(t-1)}, Z_k}$ in (29) satisfies*

$$P_{V_{k,t}|U_k^{(t-1)}, V_k^{(t-1)}, Z_k} = P_{V_{k,t}|U_k^{(t-1)}, Z_k}, \quad (46)$$

for some conditional probability measure $P_{V_{k,t}|U_k^{(t-1)}, Z_k} \in \Delta(\mathcal{V}_k|U_k^{t-1} \times \mathcal{Z}_k^{n_k})$.

Equations (45) and (46) imply, respectively, that the machine-learning algorithm and the messenger are both conforming.

Conforming algorithms are represented by probability measures conditioned on the m messages received from the server and on the local training dataset. There is no conditioning on the m messages sent to the server. The algorithms are therefore driven by the server's guidance. The same principle holds for conforming messengers. A conforming messenger is represented during communication round t by a probability measure that is conditioned on the $(t-1)$ messages received from server and the local dataset. There is no dependence on the $(t-1)$ messages transmitted. By contrast, non-conforming algorithms and messengers are represented by probability measures conditioned not only on the messages received and the local training data but also on some of messages they sent to the server. Hence, model selection reflects both the server's instructions and their own prior transmissions.

A special case of *conforming clients* is referred to as *myopic clients*.

Definition 9 (Myopic Client). *Client k , with $k \in \{1, 2, \dots, K\}$, is said to be myopic if it satisfies that*

1. the conditional probability measure $P_{\Theta_k | \mathbf{U}_k^{(m)}, \mathbf{V}_k^{(m)}, \mathbf{Z}_k}$ in (28) satisfies

$$P_{\Theta_k | \mathbf{U}_k^{(m)}, \mathbf{V}_k^{(m)}, \mathbf{Z}_k} = P_{\Theta_k | U_{k,m}, \mathbf{Z}_k}, \quad (47)$$

for some conditional probability measure $P_{\Theta_k | U_{k,m}, \mathbf{Z}_k} \in \Delta(\mathcal{M}_k | \mathcal{U}_k \times \mathcal{Z}_k^{n_k})$;
and

2. the conditional probability measure $P_{V_{k,t} | \mathbf{U}_k^{(t-1)}, \mathbf{V}_k^{(t-1)}, \mathbf{Z}_k}$ in (29) satisfies, for all $t \in \{1, 2, \dots, m\}$,

$$P_{V_{k,t} | \mathbf{U}_k^{(t-1)}, \mathbf{V}_k^{(t-1)}, \mathbf{Z}_k} = P_{V_{k,t} | U_{k,t}, \mathbf{Z}_k}, \quad (48)$$

for some conditional probability measure $P_{V_{k,t} | U_{k,t}, \mathbf{Z}_k} \in \Delta(\mathcal{V}_k | \mathcal{U}_k \times \mathcal{Z}_k^{n_k})$.

Equations (47) and (48) imply, respectively, that the machine-learning algorithm and the messenger are myopic. A myopic client is, by construction, conforming, whereas the converse does not hold in general. Myopic servers choose both the messages to be transmitted to the server and the models by sampling probability measures that are conditioned solely on the last message received from the server and the local training dataset.

The following definitions describe the main assumptions on the machine learning algorithm. The most basic classification of algorithms is between *federated algorithms* and *non-federated algorithms*.

Definition 10 (Non-federated Algorithms). *The machine learning algorithm of client k , with $k \in \{1, 2, \dots, K\}$, is said to be nonfederated if the conditional probability measure $P_{\Theta_k | \mathbf{U}_k^{(m)}, \mathbf{V}_k^{(m)}, \mathbf{Z}_k}$ in (28) satisfies*

$$P_{\Theta_k | \mathbf{U}_k^{(m)}, \mathbf{V}_k^{(m)}, \mathbf{Z}_k} = P_{\Theta_k | \mathbf{Z}_k}, \quad (49)$$

for some conditional probability measure $P_{\Theta_k | \mathbf{Z}_k} \in \Delta(\mathcal{M}_k | \mathcal{Z}_k^{n_k})$.

Non-federated learning algorithms are represented by probability measures exclusively conditioned on the local training datasets. That is, regardless of the messages sent to or received from the server. Alternatively, *federated algorithms* choose models by sampling a probability measure that is conditioned on the local training dataset and some (or all) messages sent to and received from the server.

A special case of non-federated algorithm is often referred to as Gibbs algorithms and play a central role in this work. Section 4 thoroughly describes such algorithms and their fundamental properties.

3 Generalization Error

The key feature of federated learning is that the local model chosen by client k , $\theta_k \in \mathcal{M}_k$ with $k \in \{1, 2, \dots, K\}$, is not determined solely by the client's own training dataset $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$. Instead, due to the server-mediated communication, it is fundamentally influenced by the collection of all training datasets $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K$ across all clients in the federation. See the factorization in (31). A statistical description of how clients choose their models depending on the training datasets of all clients is referred to as the *meta-federated learning algorithm* (or meta-algorithm) of client k .

3.1 Meta-Federated Learning Algorithm

The meta-federated learning algorithm of client k consists of the marginal (conditional) probability measure in $\Delta(\mathcal{M}_k | \mathcal{Z}_0)$ derived from the (conditional) probability measure $P_{\Theta, \mathbf{U}^{(m)}, \mathbf{V}^{(m)} | \mathcal{Z}_0}$ in (26), which describes the whole federated learning system. In the following, such a conditional probability measure is denoted by $P_{\Theta_k | \mathcal{Z}_0}^{(m)} \in \Delta(\mathcal{M}_k | \mathcal{Z}_0)$. The superindex (m) is for highlighting the number of communication rounds. The following lemma provides an explicit expression of this conditional probability measure.

Lemma 2. *The marginal probability measure in $\Delta(\mathcal{M}_k | \mathcal{Z}_0)$ of the probability measure $P_{\Theta, \mathbf{U}^{(m)}, \mathbf{V}^{(m)} | \mathcal{Z}_0}$ in (26), denoted $P_{\Theta_k | \mathcal{Z}_0}^{(m)} \in \Delta(\mathcal{M}_k | \mathcal{Z}_0)$, is such that for all measurable sets $\mathcal{A} \subseteq \mathcal{M}_k$ and for all $\mathbf{z}_0 = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K) \in \mathcal{Z}_0$,*

$$P_{\Theta_k | \mathcal{Z}_0 = \mathbf{z}_0}^{(m)}(\mathcal{A}) = \sum_{(\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}) \in \mathcal{U}_0^m \times \mathcal{V}_0^m} a_{\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, \mathbf{z}_0} P_{\Theta_k | \mathbf{U}_k^{(m)} = \underline{\mathbf{u}}_k^{(m)}, \mathbf{V}_k^{(m)} = \underline{\mathbf{v}}_k^{(m)}, \mathbf{Z}_k = \mathbf{z}_k}(\mathcal{A}), \quad (50)$$

where the conditional probability measure $P_{\Theta_k | \mathbf{U}_k^{(m)}, \mathbf{V}_k^{(m)}, \mathbf{Z}_k}$ is in (31). Moreover,

$$a_{\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, \mathbf{z}_0} \geq 0 \text{ and } \sum_{(\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}) \in \mathcal{U}_0^m \times \mathcal{V}_0^m} a_{\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, \mathbf{z}_0} = 1.$$

Proof: The proof is presented in Appendix B. ■

The main observation of Lemma 2 is that the *meta-algorithm* of client k , $P_{\Theta_k | \mathcal{Z}_0 = \mathbf{z}_0}^{(m)}$ in (50), is a convex combination of all possible instances of client k 's

algorithm $P_{\Theta_k|U_k^{(m)}, V_k^{(m)}, Z_k}$ in (31). More specifically, given the local training dataset \mathbf{z}_k in (1); and the sequence of messages $\underline{\mathbf{v}}^{(m)}$ and $\underline{\mathbf{u}}^{(m)}$ in (13) and (15), respectively; an instance of client k 's algorithm is the probability measure $P_{\Theta_k|U_k^{(m)}=\underline{\mathbf{u}}_k^{(m)}, V_k^{(m)}=\underline{\mathbf{v}}_k^{(m)}, Z_k=\mathbf{z}_k} \in \Delta(\mathcal{M}_k)$. Given the aggregated training dataset \mathbf{z}_0 in (2), there are at most $2^m \sum_{k=1}^K \bar{R}_k + \underline{R}_k$ possible instances, where \bar{R}_k and \underline{R}_k are the uplink and downlink information transmission rates of client k described in (11) and (12). The key observation here is that the meta-federated learning algorithm of client k depends on the training datasets of all clients via the weights $a_{\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, \mathbf{z}_0}$ with which the different instances of the algorithm $P_{\Theta_k|U_k^{(m)}, V_k^{(m)}, Z_k}$ are combined. Interestingly, such weights are determined by the client's messenger and the server. This is clearer from (31) and the following expression

$$a_{\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, \mathbf{z}_0} \triangleq \prod_{t=1}^m P_{U^{(t)}|\underline{U}^{(t-1)}=\underline{\mathbf{u}}^{(t-1)}, \underline{V}^{(t)}=\underline{\mathbf{v}}^{(t)}}(\{\mathbf{u}^{(t)}\}) \prod_{j=1}^K P_{V_{j,t}|\underline{U}_j^{(t-1)}=\underline{\mathbf{u}}_j^{(t-1)}, \underline{V}_j^{(t-1)}=\underline{\mathbf{v}}_j^{(t-1)}, Z_j=\mathbf{z}_j}(\{v_{j,t}\}), \quad (51)$$

which appears also in the proof Lemma 2, in Appendix B, Equation (123).

An additional observation from Lemma 2 is that the probability measure from which client k draws its models, *i.e.*, its meta-algorithm $P_{\Theta_k|Z_0=\mathbf{z}_0}^{(m)}$, is not arbitrary in $\Delta(\mathcal{M}_k)$. In fact, the only probability distributions (meta-algorithms) achievable by client k are mixtures of the instances of its algorithm $P_{\Theta_k|U_k^{(m)}, V_k^{(m)}, Z_k}$ in (50). Consequently, the realizable set of meta-algorithms is a subset of $\Delta(\mathcal{M}_k)$. In particular, while a meta-algorithm, induced by the aggregated training dataset, must lie in this subset, an algorithm trained using the aggregated dataset may, in principle, range over the entire set $\Delta(\mathcal{M}_k)$. More importantly, if the individual target learning algorithm is not a mixture of the instances of the local learning algorithm $P_{\Theta_k|U_k^{(m)}, V_k^{(m)}, Z_k}$ in (50), then it is not a meta-algorithm and thus not achievable by client k . This highlights an important weakness of federated learning that often goes uncommented in existing literature.

Interestingly, the definition of generalization error in federated learning can be expressed in terms of such a statistical description, as shown in the remaining of this section.

3.2 Expected Empirical Risks

In a practical federated setting, the generalization capability of a client-side learning algorithm is evaluated via *testing datasets*, which differ from the *training dataset* [57]. To formally define the expectation of both local testing and local training empirical risks, consider the following functionals, for some fixed dataset

$\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$ and some fixed model $\boldsymbol{\theta}_k \in \mathcal{M}_k$,

$$\mathbb{R}_{k, \mathbf{z}_k} : \begin{cases} \Delta(\mathcal{M}_k) \longrightarrow [0, +\infty) \\ P \longmapsto \int \mathbb{L}_k(\mathbf{z}_k, \boldsymbol{\theta}) \, dP(\boldsymbol{\theta}) \end{cases} \quad (52)$$

and

$$\mathbb{R}_{k, \boldsymbol{\theta}_k} : \begin{cases} \Delta(\mathcal{Z}_k^{n_k}) \longrightarrow [0, +\infty) \\ P \longmapsto \int \mathbb{L}_k(\mathbf{z}, \boldsymbol{\theta}_k) \, dP(\mathbf{z}), \end{cases} \quad (53)$$

where the function \mathbb{L}_k is defined in (6).

Using this notation, the expected testing and training empirical risks are defined hereunder.

Definition 11 (Expected Testing and Training Empirical Risks). *Consider the aggregated training datasets $\mathbf{z}_0 = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K) \in \mathcal{Z}_0$ and a testing dataset $\widehat{\mathbf{z}}_k \in \mathcal{Z}_k^{n_k}$, with $k \in \{1, 2, \dots, K\}$. In the federated learning system described by the measure $P_{\boldsymbol{\Theta}, \underline{\mathbf{v}}^{(m)}, \underline{\mathbf{v}}^{(m)} | \mathcal{Z}_0}$ in (26), the expected training empirical risk at client k is $\mathbb{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\Theta}_k | \mathcal{Z}_0 = \mathbf{z}_0}^{(m)} \right)$, where the probability measure $P_{\boldsymbol{\Theta}_k | \mathcal{Z}_0 = \mathbf{z}_0}^{(m)}$ is defined in (50) and the functional $\mathbb{R}_{k, \mathbf{z}_k}$ is defined in (52). Alternatively, the expected testing empirical risk at client k is $\mathbb{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\Theta}_k | \mathcal{Z}_0 = \mathbf{z}_0}^{(m)} \right)$.*

Consider the conditional probability measure $P_{\mathcal{Z}_0 | \boldsymbol{\Theta}_k}^{(m)} \in \Delta(\mathcal{Z}_0 | \mathcal{M}_k)$ to be such that for all $(\boldsymbol{\theta}_k, \mathbf{z}_0) \in \mathcal{M}_k \times \mathcal{Z}_0$,

$$\frac{dP_{\mathcal{Z}_0 | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(m)}}{dP_{\mathcal{Z}_0}}(\mathbf{z}_0) = \frac{dP_{\boldsymbol{\Theta}_k | \mathcal{Z}_0 = \mathbf{z}_0}^{(m)}}{dP_{\boldsymbol{\Theta}_k}}(\boldsymbol{\theta}_k), \quad (54)$$

where the conditional probability measure $P_{\boldsymbol{\Theta}_k | \mathcal{Z}_0}^{(m)}$ is defined in (50) and the probability measure $P_{\boldsymbol{\Theta}_k}^{(m)}$ is such that for all measurable sets $\mathcal{A} \subseteq \mathcal{M}_k$

$$P_{\boldsymbol{\Theta}_k}^{(m)}(\mathcal{A}) = \int P_{\boldsymbol{\Theta}_k | \mathcal{Z}_0 = \mathbf{z}_0}^{(m)}(\mathcal{A}) \, dP_{\mathcal{Z}_0}(\mathbf{z}_0). \quad (55)$$

The existence of the conditional probability measure $P_{\mathcal{Z}_0 | \boldsymbol{\Theta}_k}^{(m)}$ in (54) is guaranteed by [58, Theorem 11]. Moreover, consider the conditional probability measure

$$P_{\mathcal{Z}_k | \boldsymbol{\Theta}_k}^{(m)} \in \Delta(\mathcal{Z}_k^{n_k} | \mathcal{M}_k), \quad (56)$$

which is the marginal in $\Delta(\mathcal{Z}_k^{n_k} | \mathcal{M}_k)$ of the conditional probability measure $P_{\mathcal{Z}_0 | \boldsymbol{\Theta}_k}^{(m)}$ in (54).

The following lemmas, establish a connection between the functionals $\mathbb{R}_{k, \mathbf{z}_k}$ in (52) and $\mathbb{R}_{k, \boldsymbol{\theta}_k}$ in (53).

Lemma 3. Assume that for all $(\mathbf{z}_0, \boldsymbol{\theta}_k) \in (\mathcal{Z}_0 \times \mathcal{M}_k)$, the measures $P_{\mathbf{Z}_0}$ in (27); $P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ in (50); $P_{\mathbf{Z}_0 | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(m)}$ in (54); and $P_{\boldsymbol{\Theta}_k}^{(m)}$ in (55), satisfy $P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \ll P_{\boldsymbol{\Theta}_k}^{(m)}$ and $P_{\mathbf{Z}_0 | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(m)} \ll P_{\mathbf{Z}_0}$. Then, it follows that

$$\int \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0) = \int \mathbf{R}_{k, \boldsymbol{\theta}_k} \left(P_{\mathbf{Z}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(m)} \right) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k), \quad (57)$$

where the functionals $\mathbf{R}_{k, \mathbf{z}_k}$ and $\mathbf{R}_{k, \boldsymbol{\theta}_k}$ are defined in (52) and (53), respectively; and the conditional probability measure $P_{\mathbf{Z}_k | \boldsymbol{\Theta}_k}^{(m)}$ is defined in (56).

Proof: The proof is presented in Appendix C ■

Lemma 4. Assume that for all $(\mathbf{z}_0, \boldsymbol{\theta}_k) \in (\mathcal{Z}_0 \times \mathcal{M}_k)$, the measures $P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ in (50); $P_{\boldsymbol{\Theta}_k}^{(m)}$ in (55); $P_{\mathbf{Z}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(m)}$ in (56); and $P_{\mathbf{Z}_k}$ in (32), satisfy $P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \ll P_{\boldsymbol{\Theta}_k}^{(m)}$ and $P_{\mathbf{Z}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(m)} \ll P_{\mathbf{Z}_k}$. Then, it follows that

$$\begin{aligned} & \iint \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) \\ &= \int \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\Theta}_k}^{(m)} \right) dP_{\mathbf{Z}_k}(\mathbf{z}_k) \end{aligned} \quad (58)$$

$$= \int \mathbf{R}_{k, \boldsymbol{\theta}_k} \left(P_{\mathbf{Z}_k} \right) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) \quad (59)$$

$$= \iint \mathbf{R}_{k, \widehat{\boldsymbol{\theta}}_k} \left(P_{\mathbf{Z}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(m)} \right) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\widehat{\boldsymbol{\theta}}_k), \quad (60)$$

where the functionals $\mathbf{R}_{k, \mathbf{z}_k}$ and $\mathbf{R}_{k, \boldsymbol{\theta}_k}$ are defined in (52) and (53), respectively; and the probability measure $P_{\mathbf{Z}_0}$ is defined in (27).

Proof: The proof is presented in Appendix D ■

3.3 Definition of Generalization Error

The generalization error of client k , with $k \in \{1, 2, \dots, K\}$, is the expectation with respect to the joint probability measure of local models and all clients' training datasets of the difference between the local testing and local training empirical risks. In particular, for client k , let the meta-algorithm be $P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0}^{(m)} \in \Delta(\mathcal{M}_k | \mathcal{Z}_0)$, let $\widehat{\mathbf{z}}_k \in (\mathcal{X} \times \mathcal{Y})^{n_k}$ be a testing dataset, and $\mathbf{z}_k \in (\mathcal{X} \times \mathcal{Y})^{n_k}$ a training dataset, where \mathbf{z}_k and $\widehat{\mathbf{z}}_k$ are components of \mathbf{z}_0 and $\widehat{\mathbf{z}}_0$ as in (2). A typical evaluation of the generalization capability of the instance $P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ consists in the difference

$$\mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) - \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right), \quad (61)$$

where the functionals $\mathbf{R}_{k, \widehat{\mathbf{z}}_k}$ and $\mathbf{R}_{k, \mathbf{z}_k}$ are defined in (52). In a nutshell, the evaluation measures the variation of the expected empirical risk induced by $P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ when the dataset changes from \mathbf{z}_k (training) to $\widehat{\mathbf{z}}_k$ (testing). A small variation indicates good generalization for client k . While a small value of (61)

on specific testing datasets does not guarantee generalization to all datasets, this is a widely used practical criterion.

The analysis can be deepened under additional assumptions. A common one and the one adopted here to define the *generalization error* is that, for each client k , training and testing datasets are independently drawn from the same probability measure. In our federated setting, models are produced conditionally on the joint dataset $\mathbf{z}_0 \sim P_{\mathbf{Z}_0}$, while the client- k testing dataset $\widehat{\mathbf{z}}_k$ is drawn independently from the marginal $P_{\mathbf{Z}_k}$.

Definition 12 (Generalization Error). *In the federated learning system described by the measure $P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0}$ in (26), under the assumption that datasets $\mathbf{z}_0 = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K) \in \mathcal{Z}_0$ are obtained by sampling a probability measure $P_{\mathbf{Z}_0} \in \Delta(\mathcal{Z}_0)$, the generalization error at client k is*

$$\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0; P_{\mathbf{Z}_0} \right) \triangleq \iint \left(\mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\underline{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) - \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\underline{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) \right) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0), \quad (62)$$

where the functionals $\mathbf{R}_{k, \mathbf{z}_k}$ and $\mathbf{R}_{k, \widehat{\mathbf{z}}_k}$ are defined in (52); the probability measure $P_{\underline{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ is defined in (50); and the measure $P_{\mathbf{Z}_k}$ is defined in (32).

From Definition 12, it follows that the generalization error is the expectation of the difference between the expected testing empirical error $\mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\underline{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right)$ and the expected training empirical error $\mathbf{R}_{k, \mathbf{z}_k} \left(P_{\underline{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right)$ induced by the meta-algorithm $P_{\underline{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ of client k under the assumption that the training dataset $\mathbf{z}_0 = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K)$ and the testing dataset $\widehat{\mathbf{z}}_k$ are independently sampled from the measures $P_{\mathbf{Z}_0}$ and $P_{\mathbf{Z}_k}$, respectively. Note that this assumption concerns only the training and testing datasets; it does not require clients' datasets to be independent or identically distributed. This definition, while written explicitly using the meta-algorithm, is identical to the one typically used in federated learning literature, cf. [37], [44], and [45].

The generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0; P_{\mathbf{Z}_0} \right)$ in (62) can also be expressed in terms of the functional \mathbf{R}_{k, θ_k} in (53), using Lemma 3 and Lemma 4.

Lemma 5. *The generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0; P_{\mathbf{Z}_0} \right)$ in (62) satisfies*

$$\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0; P_{\mathbf{Z}_0} \right) = \iint \left(\mathbf{R}_{k, \widehat{\theta}_k} \left(P_{\mathbf{Z}_k | \underline{\Theta}_k = \theta_k}^{(m)} \right) - \mathbf{R}_{k, \theta_k} \left(P_{\mathbf{Z}_k | \underline{\Theta}_k = \theta_k}^{(m)} \right) \right) dP_{\underline{\Theta}_k}^{(m)}(\widehat{\theta}_k) dP_{\underline{\Theta}_k}^{(m)}(\theta_k), \quad (63)$$

where the functionals $R_{k, \hat{\theta}_k}$ and R_{k, θ_k} are defined in (53); the probability measure $P_{\theta_k}^{(m)} \in \Delta(\mathcal{M}_k)$ is defined in (55); and the conditional probability measure $P_{\mathbf{Z}_k | \theta_k}^{(m)} \in \Delta(\mathcal{Z}_k^{n_k} | \mathcal{M}_k)$ is defined in (56).

Proof: The proof is presented in Appendix E ■

The expression of the generalization error in (63) is, from a system-implementation perspective, less intuitive than the form in (62). Nevertheless, it is pivotal for the conclusions drawn in the following sections. A statistical intuition for (63) can be obtained as follows. Given the meta algorithm $P_{\theta_k | \mathbf{Z}_0}^{(m)}$ in (50); the probability measure $P_{\mathbf{Z}_0}$ in (27); and a model θ_k sampled from $P_{\theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ for some aggregated training dataset \mathbf{z}_0 , the probability measure $P_{\mathbf{Z}_k | \theta_k = \theta_k}^{(m)}$ can be interpreted as the posterior of $P_{\mathbf{Z}_k} \in \Delta(\mathcal{Z}_k^{n_k})$ after the observation of the model θ_k .

4 Gibbs Measures

4.1 Data-Dependent Gibbs Measures

The Gibbs conditional probability measure used in this work is parametrized by the empirical risk function L_k ; a σ -finite measure $Q_{\theta_k} \in \Delta(\mathcal{M}_k)$; 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 measure consider the following function:

$$K_{k, Q_{\theta_k}, \mathbf{z}_k} : \begin{cases} \mathbb{R} \rightarrow \mathbb{R} \\ t \mapsto \log \left(\int \exp(t L_k(\mathbf{z}_k, \theta_k)) dQ_{\theta_k}(\theta_k) \right). \end{cases} \quad (64)$$

Under the assumption that the reference measure Q_{θ_k} is a probability measure, the function $K_{k, Q_{\theta_k}, \mathbf{z}_k}$ in (64) is the cumulant generating function of the random variable $L_k(\mathbf{z}_k, \theta_k)$, for some fixed dataset $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$, when the model θ_k is sampled from Q_{θ_k} . Using this notation, the definition of Gibbs conditional probability measures on the models is presented hereunder.

Definition 13. *Given the function L_k in (6), with $k \in \{1, 2, \dots, K\}$; a σ -finite measure $Q_{\theta_k} \in \Delta(\mathcal{M}_k)$; and a $\lambda_k \in \mathbb{R} \setminus \{0\}$, the probability measure $P_{\theta_k | \mathbf{Z}_k}^{(Q_{\theta_k}, \lambda_k)} \in \Delta(\mathcal{M}_k | \mathcal{Z}_k^{n_k})$ is said to be an $(L_k, Q_{\theta_k}, \lambda_k)$ -Gibbs conditional probability measure if*

$$\forall \mathbf{z}_k \in \mathcal{S}_k, K_{k, Q_{\theta_k}, \mathbf{z}_k} \left(-\frac{1}{\lambda_k} \right) < +\infty; \quad (65)$$

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

$$\frac{dP_{\theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\theta_k}, \lambda_k)}}{dQ_{\theta_k}}(\theta_k) = \exp \left(-\frac{1}{\lambda_k} L_k(\mathbf{z}_k, \theta_k) - K_{k, Q_{\theta_k}, \mathbf{z}_k} \left(-\frac{1}{\lambda_k} \right) \right), \quad (66)$$

where the function $K_{k, Q_{\theta_k}, \mathbf{z}_k}$ is defined in (64).

Note that, while $P_{\Theta_k|Z_k}^{(Q_{\Theta_k}, \lambda_k)}$ in (66) is referred to as a Gibbs conditional probability measure, the measure $P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}$, obtained by conditioning upon a given dataset $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$, is referred to as a Gibbs probability measure.

On another note, the Gibbs probability measure $P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}$ in (66) can be related to the following optimization problem, under the assumption that $\lambda_k > 0$:

$$\min_{P \in \Delta_{Q_{\Theta_k}}(M_k)} R_{k, z_k}(P) + \lambda_k D(P \parallel Q_{\Theta_k}), \quad (67)$$

where the functional R_{k, z_k} is defined in (52). Alternatively, when $\lambda_k < 0$, such a Gibbs probability measure can be related to the following optimization problem:

$$\max_{P \in \Delta_{Q_{\Theta_k}}(M_k)} R_{k, z_k}(P) + \lambda_k D(P \parallel Q_{\Theta_k}). \quad (68)$$

The following lemma shows the connections between the Gibbs probability measure $P_{\Theta_k|Z_k}^{(Q_{\Theta_k}, \lambda_k)}$ in (66) and the optimization problems mentioned above.

Lemma 6. *Assume that the optimization problem in (67) (respectively, in (68)) admits a solution. Then, if $\lambda_k > 0$ (respectively, if $\lambda_k < 0$) the probability measure $P_{\Theta_k|Z_k}^{(Q_{\Theta_k}, \lambda_k)}$ in (66) is the unique solution.*

Proof: The proof is immediate from [35, Lemma 1]. ■

This result has also been reported in a more general context in [59, 60]. In particular, when $\lambda_k > 0$, the conditional probability measure $P_{\Theta_k|Z_k}^{(Q_{\Theta_k}, \lambda_k)}$ in (66) represents the long-run distribution of a stochastic gradient descent algorithm [61] that aims to minimize the expected empirical risk. In statistical learning, such a distribution is often referred to as the *Gibbs algorithm* [57]. Alternatively, when $\lambda_k < 0$, such a conditional measure represents an adversarial algorithm that aims to maximize the expected empirical risk.

The following lemma introduces a connection between the function K_{k, Q_{Θ_k}, z_k} in (64) and the optimization problems mentioned above.

Lemma 7. *Given an $(L_k, Q_{\Theta_k}, \lambda_k)$ -Gibbs probability measure, denoted by $P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}$, with $\mathbf{z}_k \in \Delta(\mathcal{Z}_k^{n_k})$, the following holds,*

$$-\lambda_k K_{k, Q_{\Theta_k}, z_k} \left(-\frac{1}{\lambda_k} \right) = R_{k, z_k} \left(P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) + \lambda_k D \left(P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)} \parallel Q_{\Theta_k} \right) \quad (69)$$

$$= R_{k, z_k} (Q_{\Theta_k}) - \lambda_k D \left(Q_{\Theta_k} \parallel P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)} \right), \quad (70)$$

where the function K_{k, Q_{Θ_k}, z_k} is defined in (64); and the functional R_{k, z_k} is defined in (52). Moreover, if $\lambda_k > 0$,

$$-\lambda_k K_{k, Q_{\Theta_k}, z_k} \left(-\frac{1}{\lambda_k} \right) = \min_{P \in \Delta_{Q_{\Theta_k}}(M_k)} R_{k, z_k}(P) + \lambda_k D(P \parallel Q_{\Theta_k}). \quad (71)$$

Alternatively, if $\lambda_k < 0$,

$$-\lambda_k \mathsf{K}_{k, \mathcal{Q}_{\mathbf{z}_k}, \mathbf{z}_k} \left(-\frac{1}{\lambda_k} \right) = \max_{P \in \Delta_{\mathcal{Q}_{\mathbf{z}_k}}(\mathcal{M}_k)} \mathsf{R}_{k, \mathbf{z}_k}(P) + \lambda_k D(P \parallel \mathcal{Q}_{\mathbf{z}_k}). \quad (72)$$

Proof: The proof follows immediately from [35, Lemma 2]. \blacksquare

4.2 Model-Dependent Gibbs Measures

In this section the conditional Gibbs measure is in $\Delta(\mathcal{Z}_k^{n_k} | \mathcal{M}_k)$ and is parametrized by the empirical risk function L_k ; a σ -finite measure $\mathcal{Q}_{\mathbf{z}_k} \in \Delta(\mathcal{Z}_k^{n_k})$; and a model $\boldsymbol{\theta}_k \in \mathcal{M}_k$. Using this notation, consider the following function:

$$\mathsf{K}_{k, \mathcal{Q}_{\mathbf{z}_k}, \boldsymbol{\theta}_k} : \begin{cases} \mathbb{R} \rightarrow \mathbb{R} \\ t \mapsto \log \left(\int \exp(t \mathsf{L}_k(\mathbf{z}_k, \boldsymbol{\theta}_k)) d\mathcal{Q}_{\mathbf{z}_k}(\mathbf{z}_k) \right). \end{cases} \quad (73)$$

Under the assumption that the reference measure $\mathcal{Q}_{\mathbf{z}_k}$ is a probability measure, the function $\mathsf{K}_{k, \mathcal{Q}_{\mathbf{z}_k}, \boldsymbol{\theta}_k}$ in (73) is the cumulant generating function of the random variable $\mathsf{L}_k(\mathbf{z}_k, \boldsymbol{\theta}_k)$, for some fixed model $\boldsymbol{\theta}_k \in \mathcal{M}_k$ when the dataset \mathbf{z}_k is sampled from $\mathcal{Q}_{\mathbf{z}_k}$. Using this notation, the definition of the Gibbs conditional probability measures on the datasets is presented hereunder.

Definition 14. *Given the empirical risk function L_k in (6) of client k , with $k \in \{1, 2, \dots, K\}$; a σ -finite measure $\mathcal{Q}_{\mathbf{z}_k}$ on $\Delta(\mathcal{Z}_k^{n_k})$; and a $\alpha_k \in \mathbb{R} \setminus \{0\}$, the probability measure $P_{\widehat{\mathcal{Z}}_k | \boldsymbol{\theta}_k}^{(\mathcal{Q}_{\mathbf{z}_k}, \alpha_k)} \in \Delta(\mathcal{Z}_k^{n_k} | \mathcal{M}_k)$ is said to be an $(\mathsf{L}_k, \mathcal{Q}_{\mathbf{z}_k}, \alpha_k)$ -Gibbs conditional probability measure if*

$$\forall \boldsymbol{\theta}_k \in \mathcal{B}_k, \mathsf{K}_{k, \mathcal{Q}_{\mathbf{z}_k}, \boldsymbol{\theta}_k} \left(-\frac{1}{\alpha_k} \right) < +\infty; \quad (74)$$

for some set $\mathcal{B}_k \subseteq \mathcal{M}_k$; and for all $(\mathbf{z}_k, \boldsymbol{\theta}_k) \in \text{supp } \mathcal{Q}_{\mathbf{z}_k} \times \mathcal{B}_k$,

$$\frac{dP_{\widehat{\mathcal{Z}}_k | \boldsymbol{\theta}_k}^{(\mathcal{Q}_{\mathbf{z}_k}, \alpha_k)}}{d\mathcal{Q}_{\mathbf{z}_k}}(\mathbf{z}_k) = \exp \left(-\frac{1}{\alpha_k} \mathsf{L}_k(\mathbf{z}_k, \boldsymbol{\theta}_k) - \mathsf{K}_{k, \mathcal{Q}_{\mathbf{z}_k}, \boldsymbol{\theta}_k} \left(-\frac{1}{\alpha_k} \right) \right), \quad (75)$$

where the function $\mathsf{K}_{k, \mathcal{Q}_{\mathbf{z}_k}, \boldsymbol{\theta}_k}$ is defined in (73).

Note that, while $P_{\widehat{\mathcal{Z}}_k | \boldsymbol{\theta}_k}^{(\mathcal{Q}_{\mathbf{z}_k}, \alpha_k)}$ in (75) is referred to as Gibbs conditional probability measure, the measure $P_{\widehat{\mathcal{Z}}_k | \boldsymbol{\theta}_k = \boldsymbol{\theta}_k}^{(\mathcal{Q}_{\mathbf{z}_k}, \alpha_k)}$, obtained by conditioning upon a given model $\boldsymbol{\theta}_k \in \mathcal{M}_k$, is referred to as Gibbs probability measure.

Using this notation, the Gibbs probability measures $P_{\widehat{\mathcal{Z}}_k | \boldsymbol{\theta}_k}^{(\mathcal{Q}_{\mathbf{z}_k}, \alpha_k)}$ in (75) can be related to the following optimization problems, under the assumption that $\alpha_k > 0$:

$$\min_{P \in \Delta_{\mathcal{QZ}_k}(\mathcal{Z}_k^{n_k})} \mathsf{R}_{k,\theta_k}(P) + \alpha_k D(P \parallel \mathcal{QZ}_k), \quad (76)$$

where the functional R_{k,θ_k} is defined in (53). Alternatively, when $\alpha_k < 0$, such a Gibbs probability measure can be related to the following optimization problem:

$$\max_{P \in \Delta_{\mathcal{QZ}_k}(\mathcal{Z}_k^{n_k})} \mathsf{R}_{k,\theta_k}(P) + \alpha_k D(P \parallel \mathcal{QZ}_k). \quad (77)$$

The following lemma shows the connections between the Gibbs probability measure $P_{\widehat{\mathcal{Z}}_k|\Theta_k}^{(\mathcal{QZ}_k, \alpha_k)}$ in (75) and the optimization problems mentioned above.

Lemma 8. *Assume that the optimization problem in (76) (respectively, in (77)) admits a solution. Then, if $\alpha_k > 0$ (respectively, $\alpha_k < 0$), the probability measure $P_{\widehat{\mathcal{Z}}_k|\Theta_k}^{(\mathcal{QZ}_k, \alpha_k)}$ in (75) is the unique solution.*

Proof: The proof is immediate from [35, Lemma 1]. ■

In particular, when $\alpha_k > 0$ and \mathcal{QZ}_k is a probability measure, the objective functional in (76) is the sum of two positive functionals, namely the functional R_{k,θ_k} and a relative entropy. While the former is linear, the latter is strictly convex. Hence, the objective functional is strictly convex and thus, such an optimization problem admits a unique minimizer. Alternatively, when $\alpha_k < 0$ and \mathcal{QZ}_k is a probability measure, the objective functional in (77) is strictly concave. Thus, this justifies the existence of a unique maximizer.

The following lemma introduces a connection between the function $\mathsf{K}_{k,\mathcal{QZ}_k,\theta_k}$ in (73) and the optimization problems mentioned above.

Lemma 9. *Given an $(\mathsf{L}_k, \mathcal{QZ}_k, \alpha_k)$ -Gibbs probability measure, denoted by $P_{\widehat{\mathcal{Z}}_k|\Theta_k=\theta_k}^{(\mathcal{QZ}_k, \alpha_k)}$ with $\theta_k \in \Delta(\mathcal{M}_k)$, the following holds,*

$$-\alpha_k \mathsf{K}_{k,\mathcal{QZ}_k,\theta_k} \left(-\frac{1}{\alpha_k} \right) = \mathsf{R}_{k,\theta_k} \left(P_{\widehat{\mathcal{Z}}_k|\Theta_k=\theta_k}^{(\mathcal{QZ}_k, \alpha_k)} \right) + \alpha_k D \left(P_{\widehat{\mathcal{Z}}_k|\Theta_k=\theta_k}^{(\mathcal{QZ}_k, \alpha_k)} \parallel \mathcal{QZ}_k \right) \quad (78)$$

$$= \mathsf{R}_{k,\theta_k}(\mathcal{QZ}_k) - \alpha_k D \left(\mathcal{QZ}_k \parallel P_{\widehat{\mathcal{Z}}_k|\Theta_k=\theta_k}^{(\mathcal{QZ}_k, \alpha_k)} \right), \quad (79)$$

where the function $\mathsf{K}_{k,\mathcal{QZ}_k,\theta_k}$ is defined in (73); and the functional R_{k,θ_k} is defined in (53). Moreover, if $\alpha_k > 0$,

$$-\alpha_k \mathsf{K}_{k,\mathcal{QZ}_k,\theta_k} \left(-\frac{1}{\alpha_k} \right) = \min_{P \in \Delta_{\mathcal{QZ}_k}(\mathcal{Z}_k^{n_k})} \mathsf{R}_{k,\theta_k}(P) + \alpha_k D(P \parallel \mathcal{QZ}_k). \quad (80)$$

Alternatively, if $\alpha_k < 0$,

$$-\alpha_k \mathsf{K}_{k,\mathcal{QZ}_k,\theta_k} \left(-\frac{1}{\alpha_k} \right) = \max_{P \in \Delta_{\mathcal{QZ}_k}(\mathcal{Z}_k^{n_k})} \mathsf{R}_{k,\theta_k}(P) + \alpha_k D(P \parallel \mathcal{QZ}_k). \quad (81)$$

Proof: The proof follows immediately from [35, Lemma 2]. ■

5 Method of Gaps

The method of gaps was introduced in [30] to obtain exact expressions for the generalization error of non-federated machine learning algorithms. The method consists in two steps, such steps are retained in our analysis for clarity of presentation. Interestingly, as shown hereunder, the same method applies to federated learning systems by leveraging the notion of meta-federated learning introduced above. The key observation is that the generalization error $\mathbf{G}_k \left(P_{\boldsymbol{\theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) can also be expressed in terms of gaps [30, Section IV-A].

5.1 Expected Empirical Risk Gaps

For all $k \in \{1, 2, \dots, K\}$, the functional $\mathbf{R}_{k, \mathbf{z}_k}$ in (52), for some fixed dataset $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$; and those of the variations of the functional $\mathbf{R}_{k, \boldsymbol{\theta}_k}$ in (53), for some fixed model $\boldsymbol{\theta}_k \in \mathcal{M}_k$, due to changes in their arguments, are referred to as *gaps*. These gaps are studied via the following functionals:

$$\mathbf{G}_k : \begin{cases} \mathcal{Z}_k^{n_k} \times \Delta(\mathcal{M}_k) \times \Delta(\mathcal{M}_k) \longrightarrow \mathbb{R} \\ (\mathbf{z}_k, P_1, P_2) \longmapsto \mathbf{R}_{k, \mathbf{z}_k}(P_1) - \mathbf{R}_{k, \mathbf{z}_k}(P_2) \end{cases}, \quad (82)$$

and

$$\mathbf{G}_k : \begin{cases} \mathcal{M}_k \times \Delta(\mathcal{Z}_k^{n_k}) \times \Delta(\mathcal{Z}_k^{n_k}) \longrightarrow \mathbb{R} \\ (\boldsymbol{\theta}_k, P_1, P_2) \longmapsto \mathbf{R}_{k, \boldsymbol{\theta}_k}(P_1) - \mathbf{R}_{k, \boldsymbol{\theta}_k}(P_2) \end{cases}, \quad (83)$$

where the functionals $\mathbf{R}_{k, \mathbf{z}_k}$ and $\mathbf{R}_{k, \boldsymbol{\theta}_k}$ are defined in (52) and (53), respectively. Although both functionals in (82) and (83) are denoted by \mathbf{G}_k , there is no ambiguity as they can be distinguished from their arguments.

The value $\mathbf{G}_k(\mathbf{z}_k, P_1, P_2)$, with \mathbf{G}_k in (82), represents the variation of the functional $\mathbf{R}_{k, \mathbf{z}_k}$ in (52), when its argument changes from P_2 to P_1 . Such a value is often referred to as a *data-dependent gap* as the functional \mathbf{G}_k in (82) is parametrized on a fixed dataset $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$. Conversely, the value $\mathbf{G}_k(\boldsymbol{\theta}_k, P_1, P_2)$, with \mathbf{G}_k in (83), represents the variation of the functional $\mathbf{R}_{k, \boldsymbol{\theta}_k}$ in (53), when its argument changes from P_2 to P_1 . Such a value is often referred to as a *model-dependent gap* as the functional \mathbf{G}_k in (83) is parametrized on a fixed model $\boldsymbol{\theta}_k \in \mathcal{M}_k$. This viewpoint aligns with the perspective of variations of an expectation due to changes in the probability measure presented in [35].

The following lemmas introduce a closed-form expression for the gap $\mathbf{G}_k(\mathbf{z}_k, P_{\boldsymbol{\theta}_k}^{(m)}, P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)})$ in (88) and the gap $\mathbf{G}_k(\boldsymbol{\theta}_k, P_{\mathbf{Z}_k}, P_{\mathbf{Z}_k | \boldsymbol{\theta}_k = \boldsymbol{\theta}_k}^{(m)})$ in (89).

Lemma 10. *Given an $(\mathbf{L}_k, Q_{\boldsymbol{\theta}_k}, \lambda_k)$ -Gibbs probability measure, denoted by $P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \in \Delta(\mathcal{M}_k)$, with $\mathbf{z}_k \in \Delta(\mathcal{Z}_k^{n_k})$. For all $P \in \Delta_{Q_{\boldsymbol{\theta}_k}}(\mathcal{M}_k)$,*

$$\mathbf{G}_k \left(\mathbf{z}_k, P, P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right)$$

$$= \lambda_k \left(D \left(P \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) + D \left(P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \parallel Q_{\Theta_k} \right) - D \left(P \parallel Q_{\Theta_k} \right) \right). \quad (84)$$

Proof: The proof follows immediately from [35, Lemma 3] \blacksquare

Lemma 11. *Given an $(L_k, Q_{\mathbf{Z}_k}, \alpha_k)$ -Gibbs probability measure, denoted by $P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \in \Delta(\mathcal{Z}_k^{n_k})$, with $\theta_k \in \Delta(\mathcal{M}_k)$. For all $P \in \Delta_{Q_{\mathbf{Z}_k}}(\mathcal{Z}_k^{n_k})$,*

$$\begin{aligned} & \mathbf{G}_k \left(\theta_k, P, P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \\ &= \alpha_k \left(D \left(P \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) + D \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) - D \left(P \parallel Q_{\mathbf{Z}_k} \right) \right). \end{aligned} \quad (85)$$

Proof: The proof follows immediately from [35, Lemma 3] \blacksquare

An important observation is that a gap $\mathbf{G}_k(\mathbf{z}_k, P_1, P_2)$ (or $\mathbf{G}_k(\theta_k, P_1, P_2)$) is a variation of the expectation of the functions $L_k(\mathbf{z}_k, \cdot) : \mathcal{M}_k \rightarrow [0, +\infty)$ (or $L_k(\cdot, \theta_k) : \mathcal{Z}_k^{n_k} \rightarrow [0, +\infty)$), with L_k in (6) for some fixed $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$ (or $\theta_k \in \mathcal{M}_k$) due to a change of the probability measure from P_2 to P_1 , both in $\Delta(\mathcal{M}_k)$ (or both in $\Delta(\mathcal{Z}_k^{n_k})$). From this perspective, the following lemmas introduce a closed-form expression for the gap $\mathbf{G}_k(\mathbf{z}_k, P_{\Theta_k}^{(m)}, P_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)})$ in (88) and the gap $\mathbf{G}_k(\theta_k, P_{\mathbf{Z}_k}, P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)})$ in (89).

Lemma 12. *Consider the probability measures P_1 and P_2 both in $\Delta_{Q_{\Theta_k}}(\mathcal{M}_k)$ and a dataset $\mathbf{z}_k \in \mathcal{Z}_k^{n_k}$, with $k \in \{1, 2, \dots, K\}$. Then,*

$$\begin{aligned} \mathbf{G}_k(\mathbf{z}_k, P_1, P_2) &= \lambda_k \left(D \left(P_1 \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) - D \left(P_2 \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right. \\ &\quad \left. + D(P_2 \parallel Q_{\Theta_k}) - D(P_1 \parallel Q_{\Theta_k}) \right), \end{aligned} \quad (86)$$

where $P_{\Theta_k | \mathbf{Z}_k}^{(Q_{\Theta_k}, \lambda_k)} \in \Delta(\mathcal{M}_k | \mathcal{Z}_k^{n_k})$ is the $(L_k, Q_{\Theta_k}, \lambda_k)$ -Gibbs conditional probability measure defined in (66); and the functional \mathbf{G}_k is defined in (82).

Proof: The proof is presented in Appendix F. \blacksquare

Lemma 13. *Consider the probability measures P_1 and P_2 both in $\Delta_{Q_{\mathbf{Z}_k}}(\mathcal{Z}_k^{n_k})$ and a model $\theta_k \in \mathcal{M}_k$, with $k \in \{1, 2, \dots, K\}$. Then,*

$$\begin{aligned} \mathbf{G}_k(\theta_k, P_1, P_2) &= \alpha_k \left(D \left(P_1 \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) - D \left(P_2 \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \right. \\ &\quad \left. + D(P_2 \parallel Q_{\mathbf{Z}_k}) - D(P_1 \parallel Q_{\mathbf{Z}_k}) \right), \end{aligned} \quad (87)$$

where $P_{\hat{\mathbf{Z}}_k | \Theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \in \Delta(\mathcal{Z}_k^{n_k} | \mathcal{M}_k)$ is the $(L_k, Q_{\mathbf{Z}_k}, \alpha_k)$ -Gibbs conditional probability measure defined in (75); and the functional \mathbf{G}_k is defined in (83).

Proof: The proof is presented in Appendix G. \blacksquare

5.2 Step One

The *step one* of the method of gaps consists in writing the generalization error as the expectation of a gap. The following lemmas can be assimilated to such a step.

Lemma 14. *The generalization error $\mathbf{G}_k \left(P_{\underline{\theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) satisfies*

$$\mathbf{G}_k \left(P_{\underline{\theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) = \int \mathbf{G}_k \left(\mathbf{z}_k, P_{\underline{\theta}_k}^{(m)}, P_{\underline{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0), \quad (88)$$

where the probability measure $P_{\underline{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ is defined in (50); the probability measure $P_{\underline{\theta}_k}^{(m)} \in \Delta(\mathcal{M}_k)$ is defined in (55); and the functional \mathbf{G}_k is defined in (82).

Proof: The proof is presented in Appendix H ■

Lemma 15. *The generalization error $\mathbf{G}_k \left(P_{\underline{\theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) satisfies*

$$\mathbf{G}_k \left(P_{\underline{\theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) = \int \mathbf{G}_k \left(\theta_k, P_{\mathbf{Z}_k}, P_{\mathbf{Z}_k | \underline{\theta}_k = \theta_k}^{(m)} \right) dP_{\underline{\theta}_k}^{(m)}(\theta_k), \quad (89)$$

where the probability measure $P_{\mathbf{Z}_k} \in \Delta(\mathcal{Z}_k^{n_k})$ is defined in (27); the probability measure $P_{\underline{\theta}_k}^{(m)} \in \Delta(\mathcal{M}_k)$ is defined in (55); the conditional probability measure $P_{\mathbf{Z}_k | \underline{\theta}_k}^{(m)} \in \Delta(\mathcal{Z}_k^{n_k} | \mathcal{M}_k)$ is defined in (56); and the functional \mathbf{G}_k is defined in (83).

Proof: The proof is presented in Appendix I ■

Observe that Lemma 14 and Lemma 15 express the generalization error $\mathbf{G}_k \left(P_{\underline{\theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) as the expectation of a data-dependent gap and a model-dependent gap, respectively. The significance of this observation is that these gaps admit closed-form expressions in terms of information measures namely, relative entropy.

5.3 Step Two

The *step two* of the method of gaps consists in leveraging the properties of such gaps for expressing the expectations mentioned above in terms of information measures. Using Lemma 12 and Lemma 14 leads to the following expression for the generalization error $\mathbf{G}_k \left(P_{\underline{\theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62), which can be assimilated to the *step two* of the data-dependent method of gaps.

Lemma 16. *Consider the generalization error $\mathbf{G}_k \left(P_{\underline{\theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) and assume that, for all $\mathbf{z}_0 \in \mathcal{Z}_0$, the probability measures $P_{\underline{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ in (50)*

and $P_{\Theta_k}^{(m)}$ in (55) are both absolutely continuous with respect to the σ -finite measure Q_{Θ_k} in (66). Then,

$$\begin{aligned} \mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}; P_{Z_0} \right) &= \int \lambda_k \left(D \left(P_{\Theta_k}^{(m)} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right. \\ &\left. - D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) + D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel Q_{\Theta_k} \right) - D \left(P_{\Theta_k}^{(m)} \parallel Q_{\Theta_k} \right) \right) dP_{Z_0}(z_0), \end{aligned} \quad (90)$$

where $P_{\Theta_k | Z_k}^{(Q_{\Theta_k}, \lambda_k)} \in \Delta \left(\mathcal{M}_k | \mathcal{Z}_k^{n_k} \right)$ is the $(L_k, Q_{\Theta_k}, \lambda_k)$ -Gibbs conditional probability measure defined in (66).

Proof: The proof follows directly from Lemma 12 and Lemma 14. \blacksquare

Similarly, using Lemma 13 and Lemma 15 leads to the following expression for the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}; P_{Z_0} \right)$ in (62), which can be assimilated to the *step two* of the model-dependent method of gaps.

Lemma 17. *Consider the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}; P_{Z_0} \right)$ in (62) and assume that, for all $\theta_k \in \mathcal{M}_k$, the probability measures $P_{Z_k | \Theta_k = \theta_k}^{(m)}$ in (56) and P_{Z_k} in (32), are both absolutely continuous with respect to the σ -finite measure Q_{Z_k} in (75). Then,*

$$\begin{aligned} \mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}; P_{Z_0} \right) &= \int \alpha_k \left(D \left(P_{Z_k} \parallel P_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)} \right) \right. \\ &\left. - D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)} \right) + D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel Q_{Z_k} \right) - D \left(P_{Z_k} \parallel Q_{Z_k} \right) \right) dP_{\Theta_k}^{(m)}(\theta_k), \end{aligned} \quad (91)$$

where the probability measure $P_{\Theta_k}^{(m)}$ is defined in (55); and $P_{\tilde{Z}_k | \Theta_k}^{(Q_{Z_k}, \alpha_k)} \in \Delta \left(\mathcal{Z}_k^{n_k} | \mathcal{M}_k \right)$ is the (L_k, Q_{Z_k}, α_k) -Gibbs conditional probability measure defined in (75).

Proof: The proof follows directly from Lemma 13 and Lemma 15. \blacksquare

Lemma 16 and Lemma 17 are the fundamental results from which *data-dependent* and *model-dependent* expressions are obtained for the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}; P_{Z_0} \right)$ in (62).

The following sections introduce closed-form expressions for the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}; P_{Z_0} \right)$ in (62) using the method of gaps. The expressions for the generalization error derived from (62) are referred to as *data-dependent expressions*. Those derived from (63) are referred to as *model-dependent expressions*.

6 Data-dependent Closed-Form Expressions

The exact expressions for the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) presented in this work involve an $(\mathbf{L}_k, \mathcal{Q}_{\Theta_k}, \lambda_k)$ -Gibbs conditional probability measure defined in (66).

6.1 The First Closed-Form Expression

The first closed-form expression of $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) involves a mutual information, a lautum information and a new information measure parametrized by the same parameters of the $(\mathbf{L}_k, \mathcal{Q}_{\Theta_k}, \lambda_k)$ -Gibbs probability measure $P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_{\Theta_k}, \lambda_k)}$ in (66). Such a new information measure is denoted by $J_{\mathbf{L}_k, \mathcal{Q}_{\Theta_k}, \lambda_k}$ and satisfies

$$\begin{aligned} J_{\mathbf{L}_k, \mathcal{Q}_{\Theta_k}, \lambda_k} \left(P_{\Theta_k | \mathbf{Z}_0}^{(m)}; P_{\mathbf{Z}_0} \right) &= \int \log \left(\frac{dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k)}{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_{\Theta_k}, \lambda_k)}(\theta_k)} \right) dP_{\Theta_k}^{(m)} P_{\mathbf{Z}_0}(\theta_k, \mathbf{z}_0) \\ &\quad - \int \log \left(\frac{dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k)}{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_{\Theta_k}, \lambda_k)}(\theta_k)} \right) dP_{\Theta_k | \mathbf{Z}_0}^{(m)} P_{\mathbf{Z}_0}(\theta_k, \mathbf{z}_0), \end{aligned} \quad (92)$$

where $P_{\Theta_k}^{(m)} P_{\mathbf{Z}_0} \in \Delta(\mathcal{M}_k \times \mathcal{Z}_0)$ is a product measure formed by $P_{\mathbf{Z}_0}$ in (27) and $P_{\Theta_k}^{(m)}$ in (55); $P_{\Theta_k | \mathbf{Z}_0}^{(m)} P_{\mathbf{Z}_0} \in \Delta(\mathcal{M}_k \times \mathcal{Z}_0)$ is the joint probability measure induced by the conditional measure $P_{\Theta_k | \mathbf{Z}_0}^{(m)}$ in (50) and $P_{\mathbf{Z}_0}$; and $P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_{\Theta_k}, \lambda_k)}$ is the $(\mathbf{L}_k, \mathcal{Q}_{\Theta_k}, \lambda_k)$ -Gibbs probability measure in (66).

The following theorem consolidates the hypothesis testing nature of the new information measure introduced in (92).

Theorem 18. *Consider a free parameter $\gamma > 0$; the following constants*

$$\underline{\gamma} \triangleq \min_{(\theta_k, \mathbf{z}_0) \in \mathcal{M}_k \times \mathcal{Z}_0} \log \left(\frac{dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k)}{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_{\Theta_k}, \lambda_k)}(\theta_k)} \right) \text{ and} \quad (93)$$

$$\bar{\gamma} \triangleq \max_{(\theta_k, \mathbf{z}_0) \in \mathcal{M}_k \times \mathcal{Z}_0} \log \left(\frac{dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k)}{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_{\Theta_k}, \lambda_k)}(\theta_k)} \right); \quad (94)$$

and the set $\mathcal{A}_\gamma = \left\{ (\theta_k, \mathbf{z}_0) \in \mathcal{M}_k \times \mathcal{Z}_0 : \log \left(\frac{dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k)}{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_{\Theta_k}, \lambda_k)}(\theta_k)} \right) \geq \gamma \right\}$, where

the probability measures $P_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ is defined in (50) and $P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(\mathcal{Q}_{\Theta_k}, \lambda_k)}$ is the $(\mathbf{L}_k, \mathcal{Q}_{\Theta_k}, \lambda_k)$ -Gibbs conditional probability measure defined in (66).

Then, the value $J_{\mathbf{L}_k, \mathcal{Q}_{\Theta_k}, \lambda_k} \left(P_{\Theta_k | \mathbf{Z}_0}^{(m)}; P_{\mathbf{Z}_0} \right)$ in (92), satisfies

$$(\bar{\gamma} - \gamma) P_{\Theta_k}^{(m)} P_{\mathbf{Z}_0}(\mathcal{A}_\gamma) + (\gamma - \underline{\gamma}) P_{\Theta_k | \mathbf{Z}_0}^{(m)} P_{\mathbf{Z}_0}(\mathcal{A}_\gamma^c)$$

$$\begin{aligned} &\geq J_{L_k, Q_{\Theta_k}, \lambda_k} \left(P_{\Theta_k|Z_0}^{(m)}; P_{Z_0} \right) \geq \\ &\left(\underline{\gamma} - \gamma \right) P_{\Theta_k}^{(m)} P_{Z_0} \left(\mathcal{A}_{\underline{\gamma}} \right) + \left(\gamma - \bar{\gamma} \right) P_{\Theta_k|Z_0}^{(m)} P_{Z_0} \left(\mathcal{A}_{\bar{\gamma}} \right), \end{aligned} \quad (95)$$

where $P_{\Theta_k}^{(m)} P_{Z_0} \in \Delta(\mathcal{M}_k \times \mathcal{Z}_0)$ is a product measure formed by P_{Z_0} in (27) and $P_{\Theta_k}^{(m)}$ in (55); and $P_{\Theta_k|Z_0}^{(m)} P_{Z_0} \in \Delta(\mathcal{M}_k \times \mathcal{Z}_0)$ is the joint probability measure induced by the conditional measure $P_{\Theta_k|Z_0}^{(m)}$ in (50) and P_{Z_0} .

Proof: The proof is presented in Appendix J. ■

Using mutual information, lautum information and this new information measure, the first data-dependent closed-form expression of the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0}; P_{Z_0} \right)$ in (62) is introduced hereunder.

Theorem 19. Consider the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0}; P_{Z_0} \right)$ in (62) and assume that for all $z_0 \in \mathcal{Z}_0$, the measures $P_{\Theta_k|Z_0=z_0}^{(m)}$ in (50); $P_{\Theta_k}^{(m)}$ in (55); and Q_{Θ_k} in (66), satisfy that $P_{\Theta_k}^{(m)} \ll Q_{\Theta_k} \ll P_{\Theta_k}^{(m)} \ll P_{\Theta_k|Z_0=z_0}^{(m)} \ll Q_{\Theta_k}$. Then,

$$\begin{aligned} &\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0}; P_{Z_0} \right) \\ &= \lambda_k \left(I \left(P_{\Theta_k|Z_0}^{(m)}; P_{Z_0} \right) + L \left(P_{\Theta_k|Z_0}^{(m)}; P_{Z_0} \right) + J_{L_k, Q_{\Theta_k}, \lambda_k} \left(P_{\Theta_k|Z_0}^{(m)}; P_{Z_0} \right) \right), \end{aligned} \quad (96)$$

where $J_{L_k, Q_{\Theta_k}, \lambda_k} \left(P_{\Theta_k|Z_0}^{(m)}; P_{Z_0} \right)$ is defined in (92).

Proof: The proof is presented in Appendix K. ■

In most practical settings, the absolute-continuity assumptions in Theorem 19 are satisfied; exceptions typically arise only in specialized academic constructions [62]. Together, Theorem 18 and Theorem 19 lead to a new bound on the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0}; P_{Z_0} \right)$ in (62), which beyond its tightness, reveals the entanglement between the analysis of the generalization error and the study of mismatched hypothesis tests.

6.2 The Second Closed-Form Expression

The second data-dependent exact expression for the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0}; P_{Z_0} \right)$ in (62) involves exclusively relative entropies.

Theorem 20. Consider the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0}; P_{Z_0} \right)$ in (62) and assume that for all $z_0 \in \mathcal{Z}_0$, $P_{\Theta_k|Z_0=z_0}^{(m)} \ll Q_{\Theta_k}$, with $P_{\Theta_k|Z_0=z_0}^{(m)}$ in (50) and Q_{Θ_k} in (66). Then,

$$\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0}; P_{Z_0} \right) = \lambda_k \iint \left(D \left(P_{\Theta_k|Z_0=z_0}^{(m)} \parallel P_{\Theta_k|Z_k=\bar{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right)$$

$$-D\left(P_{\Theta_k|Z_0=z_0}^{(m)} \parallel P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}\right) dP_{Z_k}(\widehat{z}_k) dP_{Z_0}(z_0), \quad (97)$$

where $P_{\Theta_k|Z_k}^{(Q_{\Theta_k}, \lambda_k)} \in \Delta(\mathcal{M}_k | \mathcal{Z}_k^{n_k})$ is the $(L_k, Q_{\Theta_k}, \lambda_k)$ -Gibbs conditional probability measure defined in (66).

Proof: The proof is presented in Appendix L. ■

6.3 A Pythagorean Identity of Generalization Error

The following theorem, together with the converse of the Pythagorean theorem [63, Book I, Proposition 48] leads to a Pythagorean identity, shown in Figure 3, which is reminiscent of the information projections studied in [35, 64, 65].

Theorem 21. *Consider the generalization error $\mathbf{G}_k(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{Y}^{(m)}|Z_0}; P_{Z_0})$ in (62). Assume that the conditional measure $P_{\Theta_k|Z_0}^{(m)}$ in (50); the measure $P_{\Theta_k}^{(m)}$ in (55); and σ -finite measure Q_{Θ_k} in (66), satisfy for all $z_0 \in \mathcal{Z}_0$, $P_{\Theta_k|Z_0=z_0}^{(m)} \ll P_{\Theta_k}^{(m)} \ll Q_{\Theta_k}$. Then,*

$$\begin{aligned} & \frac{1}{\lambda_k} \mathbf{G}_k(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{Y}^{(m)}|Z_0}; P_{Z_0}) + \int D\left(P_{\Theta_k|Z_0=z_0}^{(m)} \parallel P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}\right) dP_{Z_0}(z_0) \\ &= \int D\left(P_{\Theta_k}^{(m)} \parallel P_{\Theta_k|Z_k=\widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)}\right) dP_{Z_k}(\widehat{z}_k) + \int D\left(P_{\Theta_k|Z_0=z_0}^{(m)} \parallel P_{\Theta_k}^{(m)}\right) dP_{Z_0}(z_0), \quad (98) \end{aligned}$$

where the measure P_{Z_k} is defined in (32) and $P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}$ is the $(L_k, Q_{\Theta_k}, \lambda_k)$ -Gibbs probability measure in (66).

Proof: The proof is presented in Appendix M. ■

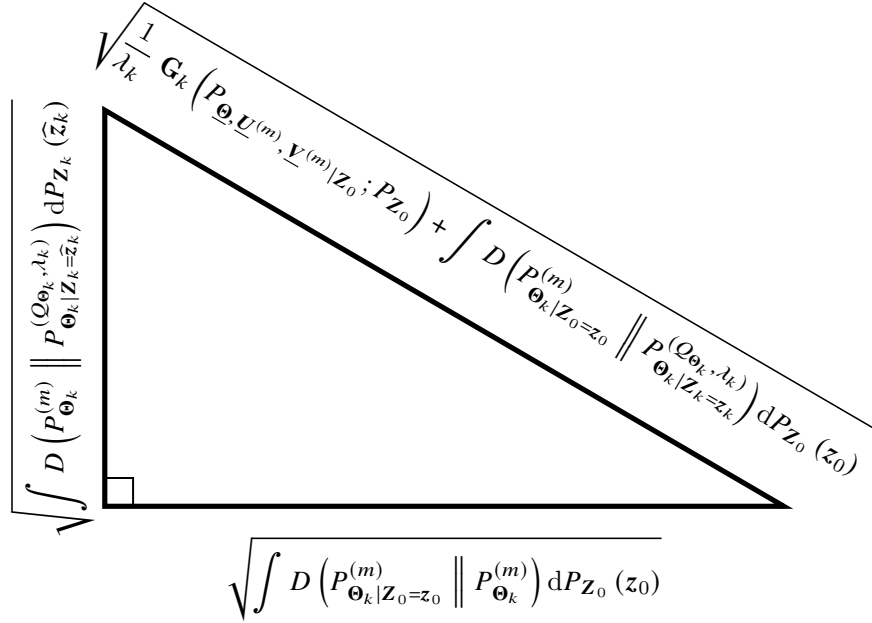


Figure 3: Geometric interpretation of Theorem 21

7 Model-dependent Closed-Form Expressions

7.1 The First Closed-Form Expression

The first closed-form expression of $\mathbf{G}_k(P_{\Theta, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0; P_{Z_0})$ in (62) involves a mutual information, a lautum information and a new information measure parametrized by the same parameters of the (L_k, Q_{Z_k}, α_k) -Gibbs probability measure $P_{\tilde{Z}_k|\Theta_k=\theta_k}^{(Q_{Z_k}, \alpha_k)}$ in (75). Such a new information measure is denoted by $J_{L_k, Q_{Z_k}, \alpha_k}$ and satisfies

$$\begin{aligned}
 J_{L_k, Q_{Z_k}, \alpha_k}(P_{Z_k|\Theta_k}^{(m)}; P_{\Theta_k}^{(m)}) &= \int \log \left(\frac{dP_{Z_k|\Theta_k=\theta_k}^{(m)}}{dP_{\tilde{Z}_k|\Theta_k=\theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) dP_{Z_k} P_{\Theta_k}^{(m)}(z_k, \theta_k) \\
 &\quad - \int \log \left(\frac{dP_{Z_k|\Theta_k=\theta_k}^{(m)}}{dP_{\tilde{Z}_k|\Theta_k=\theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) dP_{Z_k|\Theta_k}^{(m)} P_{\Theta_k}^{(m)}(z_k, \theta_k), \quad (99)
 \end{aligned}$$

where $P_{\mathbf{Z}_k}$ is the marginal probability measure in $\Delta(\mathcal{Z}_k^{n_k})$ of $P_{\mathbf{Z}_0}$ in (27), which satisfies

$$P_{\mathbf{Z}_0}(\mathcal{A}) = \int P_{\mathbf{Z}_0|\Theta_k=\theta_k}^{(m)}(\mathcal{A}) dP_{\Theta_k}^{(m)}(\theta_k), \quad (100)$$

with $P_{\mathbf{Z}_0|\Theta_k=\theta_k}^{(m)}$ defined in (54); $P_{\mathbf{Z}_k}P_{\Theta_k}^{(m)} \in \Delta(\mathcal{Z}_k^{n_k} \times \mathcal{M}_k)$ is a product measure formed by $P_{\mathbf{Z}_k}$ and $P_{\Theta_k}^{(m)}$ in (55); and $P_{\mathbf{Z}_k|\Theta_k}^{(m)}P_{\Theta_k}^{(m)} \in \Delta(\mathcal{Z}_k^{n_k} \times \mathcal{M}_k)$ is the joint probability measure induced by the conditional measure $P_{\mathbf{Z}_k|\Theta_k}^{(m)}$ in (56) and $P_{\Theta_k}^{(m)}$; and $P_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}$ is the $(\mathbb{L}_k, Q_{\mathbf{Z}_k}, \alpha_k)$ -Gibbs probability measure in (75). Note that the equality in (100) is a consequence of the equality in (54).

Using this notation the following Theorem is introduced and follows along the same lines as Theorem (18).

Theorem 22. Consider a free parameter $\gamma > 0$; the following constants

$$\underline{\gamma} \triangleq \min_{(\mathbf{z}_k, \theta_k) \in \mathcal{Z}_k^{n_k} \times \mathcal{M}_k} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}(\mathbf{z}_k)} \right) \text{ and} \quad (101)$$

$$\bar{\gamma} \triangleq \max_{(\mathbf{z}_k, \theta_k) \in \mathcal{Z}_k^{n_k} \times \mathcal{M}_k} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}(\mathbf{z}_k)} \right); \quad (102)$$

and the set $\mathcal{A}_\gamma = \left\{ (\mathbf{z}_k, \theta_k) \in \mathcal{Z}_k^{n_k} \times \mathcal{M}_k : \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}(\mathbf{z}_k)} \right) \geq \gamma \right\}$, where the probability measure $P_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}$ defined in (56) and $P_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}$ is the $(\mathbb{L}_k, Q_{\mathbf{Z}_k}, \alpha_k)$ -Gibbs conditional probability measure defined in (75). Then, the value $J_{\mathbb{L}_k, Q_{\mathbf{Z}_k}, \alpha_k}(P_{\mathbf{Z}_k|\Theta_k}^{(m)}; P_{\Theta_k}^{(m)})$ in (99), satisfies

$$\begin{aligned} & (\bar{\gamma} - \gamma) P_{\mathbf{Z}_k}P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma) + (\gamma - \underline{\gamma}) P_{\mathbf{Z}_k|\Theta_k}^{(m)}P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma^c) \\ & \geq J_{\mathbb{L}_k, Q_{\mathbf{Z}_k}, \alpha_k}(P_{\mathbf{Z}_k|\Theta_k}^{(m)}; P_{\Theta_k}^{(m)}) \geq \\ & (\underline{\gamma} - \gamma) P_{\mathbf{Z}_k}P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma^c) + (\gamma - \bar{\gamma}) P_{\mathbf{Z}_k|\Theta_k}^{(m)}P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma), \end{aligned} \quad (103)$$

where $P_{\mathbf{Z}_k}P_{\Theta_k}^{(m)} \in \Delta(\mathcal{Z}_k^{n_k} \times \mathcal{M}_k)$ is a product measure formed by $P_{\mathbf{Z}_k}$ in (32) and $P_{\Theta_k}^{(m)}$ in (55); and $P_{\mathbf{Z}_k|\Theta_k}^{(m)}P_{\Theta_k}^{(m)} \in \Delta(\mathcal{Z}_k^{n_k} \times \mathcal{M}_k)$ is the joint probability measure induced by the conditional measure $P_{\mathbf{Z}_k|\Theta_k}^{(m)}$ in (56) and $P_{\Theta_k}^{(m)}$.

Proof: The proof is presented in Appendix N. \blacksquare

Using mutual information, lautum information and this new information measure, the first model-dependent closed-form expression of the generalization error $\mathbf{G}_k(P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)}|\mathbf{Z}_0}; P_{\mathbf{Z}_0})$ in (62) is introduced below.

Theorem 23. Consider the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) and assume that for all $\theta_k \in \mathcal{M}_k$, the measures $P_{\mathbf{Z}_k}$ in (32); $P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}$ in (56); and $Q_{\mathbf{Z}_k}$ in (75); satisfy that $P_{\mathbf{Z}_k} \ll Q_{\mathbf{Z}_k} \ll P_{\mathbf{Z}_k} \ll P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \ll Q_{\mathbf{Z}_k}$. Then,

$$\begin{aligned} & \mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) \\ &= \alpha_k \left(I \left(P_{\mathbf{Z}_k | \Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) + L \left(P_{\mathbf{Z}_k | \Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) + J_{L_k, Q_{\mathbf{Z}_k}, \alpha_k} \left(P_{\mathbf{Z}_k | \Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) \right), \end{aligned} \quad (104)$$

where $J_{L_k, Q_{\mathbf{Z}_k}, \alpha_k} \left(P_{\mathbf{Z}_k | \Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right)$ is defined in (99).

Proof: The proof is presented in Appendix O. ■

Together, Theorem 22 and Theorem 23 lead to a new bound on the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62).

7.2 The Second Closed-Form Expression

The second model-dependent exact expression for the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) involves exclusively relative entropies.

Theorem 24. Consider the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) and assume that, for all $\theta_k \in \mathcal{M}_k$, the measures $P_{\mathbf{Z}_k}$ in (32); $P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}$ in (56); and $Q_{\mathbf{Z}_k}$ in (75); satisfy that $P_{\mathbf{Z}_k} \ll Q_{\mathbf{Z}_k}$ and $P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \ll Q_{\mathbf{Z}_k}$. Then,

$$\begin{aligned} \mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) &= \alpha_k \iint D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\widehat{\mathbf{Z}}_k | \Theta_k = \widehat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \\ &- D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\widehat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)} \left(\widehat{\theta}_k \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right), \end{aligned} \quad (105)$$

where the probability measures $P_{\Theta_k}^{(m)}$ is defined in (55) and $P_{\widehat{\mathbf{Z}}_k | \Theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \in \Delta \left(\mathcal{Z}_k^{n_k} | \mathcal{M}_k \right)$ is an $(L_k, Q_{\mathbf{Z}_k}, \alpha_k)$ -Gibbs conditional probability measure defined in (75).

Proof: The proof is presented in Appendix P. ■

7.3 A Pythagorean Identity of Generalization Error

The following theorem, together with the converse of the Pythagorean theorem [63, Book I, Proposition 48] leads to a Pythagorean identity, shown in Figure 4.

Theorem 25. Consider the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{Y}^{(m)} | Z_0}; P_{Z_0} \right)$ in (62) and assume that, for all $\theta_k \in \mathcal{M}_k$, the measures P_{Z_k} in (32); $P_{Z_k | \Theta_k = \theta_k}^{(m)}$ in (56); and Q_{Z_k} in (75); satisfy that $P_{Z_k | \Theta_k = \theta_k}^{(m)} \ll P_{Z_k} \ll Q_{Z_k}$. Then,

$$\begin{aligned} & \frac{1}{\alpha_k} \mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{Y}^{(m)} | Z_0}; P_{Z_0} \right) + \int D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\widehat{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)} (\theta_k) \\ &= \int D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel P_{Z_k} \right) dP_{\Theta_k}^{(m)} (\theta_k) + \int D \left(P_{Z_k} \parallel P_{\widehat{Z}_k | \Theta_k = \widehat{\theta}_k}^{(Q_{Z_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)} (\widehat{\theta}_k), \quad (106) \end{aligned}$$

where the measure $P_{\Theta_k}^{(m)}$ is defined in (55) and $P_{\widehat{Z}_k | \Theta_k = \widehat{\theta}_k}^{(Q_{Z_k}, \alpha_k)}$ is the (L_k, Q_{Z_k}, α_k) -Gibbs probability measure in (75).

Proof: The proof is presented in Appendix Q ■

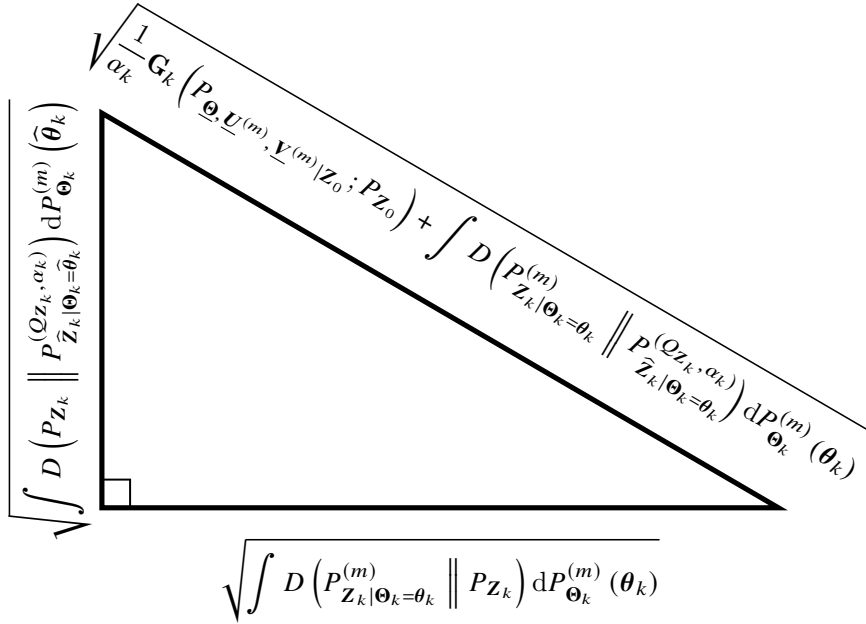


Figure 4: Geometric interpretation of Theorem 25

References

- [1] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Proceedings of the 20th International Conference on Artificial Intelligence*

- and Statistics (AISTATS)*, ser. Proceedings of Machine Learning Research, vol. 54. PMLR, Apr. 2017, pp. 1273–1282.
- [2] 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. Ozgür, 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.
- [3] J. Wang, Z. Charles, Z. Xu, G. Joshi, H. B. McMahan, B. Agüera y Arcas, M. Al-Shedivat, G. Andrew, S. Avestimehr, K. Daly, D. Data, S. Diggavi, H. Eichner, A. Gadhikar, Z. Garrett, A. M. Girgis, F. Hanzely, A. Hard, C. He, S. Horvath, Z. Huo, A. Ingerman, M. Jaggi, T. Javidi, P. Kairouz, S. Kale, S. P. Karimireddy, J. Konečný, S. Koyejo, T. Li, L. Liu, M. Mohri, H. Qi, S. J. Reddi, P. Richtárik, K. Singhal, V. Smith, M. Soltanolkotabi, W. Song, A. T. Suresh, S. U. Stich, A. Talwalkar, H. Wang, B. Woodworth, S. Wu, F. X. Yu, H. Yuan, M. Zaheer, M. Zhang, T. Zhang, C. Zheng, C. Zhu, and W. Zhu, *A Field Guide to Federated Optimization*, 1st ed. Ithaca, NY, USA: arXiv, 2021.
- [4] M. Mohri, G. Sivek, and A. T. Suresh, “Agnostic federated learning,” in *Proceedings of the International Conference on Machine Learning (ICML)*, ser. Proceedings of Machine Learning Research, vol. 97. PMLR, 2019, pp. 4615–4625.
- [5] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated optimization in heterogeneous networks,” *Proceedings of Machine Learning and Systems*, vol. 2, pp. 429–450, Mar. 2020.
- [6] S. P. Karimireddy, S. Kale, M. Mohri, S. J. Reddi, S. U. Stich, and A. T. Suresh, “SCAFFOLD: Stochastic controlled averaging for federated learning,” in *Proceedings of the International Conference on Machine Learning (ICML)*. PMLR, 2020.
- [7] A. Khaled, K. Mishchenko, and P. Richtárik, “Tighter theory for local SGD on identical and heterogeneous data,” in *International Conference on Artificial Intelligence and Statistics (AISTATS)*, vol. 108. PMLR, 2020.
- [8] J. Konečný, H. B. McMahan, D. Ramage, and P. Richtárik, “Federated optimization: Distributed machine learning for on-device intelligence,” *arXiv*, Oct. 2016.

-
- [9] S. U. Stich, “Local sgd converges fast and communicates little,” *International Conference on Learning Representations (ICLR)*, May 2019.
- [10] E. Diao, J. Ding, and V. Tarokh, “HeteroFL: Computation and communication efficient federated learning for heterogeneous clients,” in *International Conference on Learning Representations (ICLR)*, 2021.
- [11] V. Vapnik, *Estimation of dependences based on empirical data*. Springer Science & Business Media, 2006.
- [12] N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler, “Scale-sensitive dimensions, uniform convergence, and learnability,” *Journal of the ACM (JACM)*, vol. 44, no. 4, pp. 615–631, 1997.
- [13] M. Mohri and A. Rostamizadeh, “Rademacher complexity bounds for non-iid processes,” in *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 21, 2008.
- [14] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*, 1st ed. New York, NY, USA: Cambridge University Press, 2014.
- [15] V. Cherkassky, X. Shao, F. M. Mulier, and V. N. Vapnik, “Model complexity control for regression using VC generalization bounds,” *IEEE Transactions on Neural Networks*, vol. 10, no. 5, pp. 1075–1089, Sep. 1999.
- [16] B. Wei, J. Li, Y. Liu, and W. Wang, “Non-iid federated learning with sharper risk bound,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 35, no. 5, pp. 6906–6917, 2024.
- [17] X. Hu, S. Li, and Y. Liu, “Generalization bounds for federated learning: Fast rates, unparticipating clients and unbounded losses,” in *International Conference on Learning Representations (ICLR)*, 2023.
- [18] O. Bousquet and A. Elisseeff, “Stability and generalization,” *The Journal of Machine Learning Research*, vol. 2, no. 1, pp. 499–526, Mar. 2002.
- [19] M. Hardt, B. Recht, and Y. Singer, “Train faster, generalize better: Stability of stochastic gradient descent,” in *Proceedings of the 33rd International Conference on Machine Learning*, vol. 48, 2016, pp. 1225–1234.
- [20] 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.
- [21] G. K. Dziugaite and D. M. Roy, “Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data,” in *Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence*, 2017.

- [22] A. Xu and M. Raginsky, “Information-theoretic analysis of generalization capability of learning algorithms,” in *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 30, Dec. 2017, pp. 1–10.
- [23] Y. Bu, S. Zou, and V. V. Veeravalli, “Tightening mutual information-based bounds on generalization error,” *IEEE Journal on Selected Areas in Information Theory*, vol. 1, no. 1, pp. 121–130, Jan. 2020.
- [24] J. Negrea, M. Haghifam, G. K. Dziugaite, A. Khisti, and D. M. Roy, “Information-theoretic generalization bounds for SGLD via data-dependent estimates,” in *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 32, 2019, pp. 11 013–11 023.
- [25] D. Russo and J. Zou, “Controlling bias in adaptive data analysis using information theory,” in *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, vol. 51, Cadiz, Spain, May 2016, pp. 1232–1240.
- [26] G. Aminian, Y. Bu, L. Toni, M. Rodrigues, and G. Wornell, “An exact characterization of the generalization error for the Gibbs algorithm,” in *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS)*, vol. 34, Dec. 2021, pp. 8106–8118.
- [27] G. Aminian, Y. Bu, L. Toni, M. R. D. Rodrigues, and G. W. Wornell, “Information-theoretic characterizations of generalization error for the Gibbs algorithm,” *IEEE Transactions on Information Theory*, vol. 70, no. 1, pp. 632–655, 2024.
- [28] X. Zou, S. M. Perlaza, I. Esnaola, and E. Altman, “Generalization analysis of machine learning algorithms via the worst-case data-generating probability measure,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, Vancouver, Canada, Feb. 2024.
- [29] 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.
- [30] S. M. Perlaza and X. Zou, “The generalization error of machine learning algorithms,” *arXiv preprint arXiv:2411.12030*, 2024.
- [31] P. Boroumand and A. G. i Fàbregas, “Mismatched binary hypothesis testing: Error exponent sensitivity,” *IEEE Transactions on Information Theory*, vol. 68, no. 10, pp. 6738–6761, 2022.
- [32] C. E. Shannon, “A mathematical theory of communication,” *The Bell System Technical Journal*, vol. 27, pp. 379–423, Jul. 1948.
- [33] —, “A mathematical theory of communication,” *The Bell System Technical Journal*, vol. 27, pp. 623–656, Oct. 1948.

-
- [34] D. P. Palomar and S. Verdú, “Lautum information,” *IEEE Transactions on Information Theory*, vol. 54, no. 3, pp. 964–975, Mar. 2008.
- [35] 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.
- [36] I. Csiszár, “ i -divergence geometry of probability distributions and minimization problems,” *The Annals of Probability*, vol. 3, no. 1, pp. 146–158, 1975.
- [37] S. Yagli, A. Dytso, and H. V. Poor, “Information-theoretic bounds on the generalization error and privacy leakage in federated learning,” in *Proceedings of the International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*. IEEE, May 2020, pp. 1–5.
- [38] R. Pathak and M. J. Wainwright, “Fedsplit: An algorithmic framework for fast federated optimization,” in *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [39] M. Sefidgaran, R. Chor, A. Zaidi, and Y. Wan, “Lessons from generalization error analysis of federated learning: You may communicate less often!” in *Proceedings of the International Conference on Machine Learning (ICML)*, vol. 235. PMLR, Jul. 2024, pp. 44 093–44 135.
- [40] T. Wang, J. Rausch, C. Zhang, R. Jia, and D. Song, “A principled approach to data valuation for federated learning,” in *Federated Learning: Privacy and Incentive*, ser. Lecture Notes in Computer Science, Q. Yang, L. Fan, and H. Yu, Eds. Springer, 2020, vol. 12500, pp. 153–167.
- [41] J. Wang, V. Tantia, N. Ballas, and M. Rabbat, “Slowmo: Improving communication-efficient distributed SGD with slow momentum,” in *International Conference on Learning Representations (ICLR)*, 2020.
- [42] H. Yu, R. Jin, and S. Yang, “On the linear speedup analysis of communication efficient momentum SGD for distributed non-convex optimization,” in *Proceedings of the International Conference on Machine Learning (ICML)*, vol. 97. PMLR, 2019.
- [43] S. Horváth and P. Richtárik, “A better alternative to error feedback for communication-efficient distributed learning,” *arXiv*, Jun. 2020.
- [44] Z. Sun, X. Niu, and E. Wei, “Understanding generalization of federated learning via stability: Heterogeneity matters,” in *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, vol. 238. PMLR, May 2024, pp. 676–684.
- [45] L. P. Barnes, A. Dytso, and H. V. Poor, “Improved information theoretic generalization bounds for distributed and federated learning,” *arXiv*, Feb. 2022.

- [46] A. Ramezani-Kebrya, F. Liu, T. Pethick, G. Chrysos, and V. Cevher, “Federated learning under covariate shifts with generalization guarantees,” *Transactions on Machine Learning Research*, Jun. 2023.
- [47] R. Zhang, Q. Xu, J. Yao, Y. Zhang, Q. Tian, and Y. Wang, “Federated domain generalization with generalization adjustment,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Jun. 2023, pp. 3954–3963.
- [48] M. Sefidgaran, R. Chor, and A. Zaidi, “Rate-distortion theoretic bounds on generalization error for distributed learning,” in *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- [49] F. Chen, M. Luo, Z. Dong, Z. Li, and X. He, “Federated meta-learning with fast convergence and efficient communication,” *arXiv*, Feb. 2019.
- [50] H. Zhang, C. Li, N. Kan, Z. Zheng, W. Dai, J. Zou, and H. Xiong, “Improving generalization in federated learning with model-data mutual information regularization: A posterior inference approach,” in *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS)*, 2024.
- [51] H. Yuan, W. Morningstar, L. Ning, and K. Singhal, “What do we mean by generalization in federated learning?” *arXiv*, Oct. 2022.
- [52] Z. Wang, C. Long, and Y. Mao, “Generalization in federated learning: A conditional mutual information framework,” *arXiv*, Mar. 2025.
- [53] W. Liu, G. Yu, L. Wang, and R. Liao, “An information-theoretic framework for out-of-distribution generalization with applications to stochastic gradient Langevin dynamics,” *IEEE Transactions on Information Theory*, vol. 71, no. 11, pp. 8798–8817, 2025.
- [54] S. J. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan, “Adaptive federated optimization,” *International Conference on Learning Representations*, May 2021.
- [55] B. Ying, Z. Li, and H. Yang, “Exact and linear convergence for federated learning under arbitrary client participation is attainable,” Jun. 2025.
- [56] Y. J. Cho, P. Sharma, G. Joshi, Z. Xu, S. Kale, and T. Zhang, “On the convergence of federated averaging with cyclic client participation,” in *Proceedings of the 40th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, Eds., vol. 202. PMLR, 23–29 Jul 2023, pp. 5677–5721.

- [57] 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 *International Symposium on Information Theory (ISIT)*, June 2023, pp. 328–333.
- [58] 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.
- [59] 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 *2024 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2024, pp. 759–764.
- [60] —, “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.
- [61] 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*, July 2024, pp. 2168 – 2229.
- [62] Y. Polyanskiy and Y. Wu, *Information Theory*, 1st ed. Cambridge, UK: Cambridge University Press, 2023.
- [63] Euclid, *The Thirteen Books of Euclid’s Elements*, 2nd ed., T. L. Heath, Ed. Dover Publications, Inc., 1956.
- [64] N. N. Chentsov, “Nonsymmetrical distance between probability distributions, entropy and the theorem of pythagoras,” *Mathematical notes of the Academy of Sciences of the USSR*, vol. 4, no. 3, pp. 686–691, 1968.
- [65] I. Csiszár and F. Matus, “Information projections revisited,” *IEEE Transactions on Information Theory*, vol. 49, no. 6, pp. 1474–1490, 2003.
- [66] T. Bayes, “An essay towards solving a problem in the doctrine of chances,” *Biometrika*, vol. 45, no. 3-4, pp. 296–315, 1958.
- [67] R. B. Ash and C. A. Doleans-Dade, *Probability and Measure Theory*, 2nd ed. Burlington, MA, USA: Academic Press, 2000.

A Proof of Lemma 1

From Bayes' theorem [66], the joint probability measure $P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}, \mathbf{Z}_0}$ in (24) can be factorized as,

$$P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0} = P_{\underline{\Theta} | \underline{U}^{(m)}, \underline{V}^{(m)}, \mathbf{Z}_0} P_{\underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0} \quad (107a)$$

with,

$$P_{\underline{\Theta} | \underline{U}^{(m)}, \underline{V}^{(m)}, \mathbf{Z}_0} \in \Delta(\mathcal{M}_0 | \mathcal{U}_0^m \times \mathcal{V}_0^m \times \mathcal{Z}_0), \text{ and} \quad (107b)$$

$$P_{\underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0} \in \Delta(\mathcal{U}_0^m \times \mathcal{V}_0^m | \mathcal{Z}_0). \quad (107c)$$

Furthermore, under Assumption 1, the conditional probability measure $P_{\underline{\Theta} | \underline{U}^{(m)}, \underline{V}^{(m)}, \mathbf{Z}_0}$ in (107b) can be factorized as

$$P_{\underline{\Theta} | \underline{U}^{(m)}, \underline{V}^{(m)}, \mathbf{Z}_0} = \prod_{k=1}^K P_{\Theta_k | \mathbf{U}_k^{(m)}, \mathbf{V}_k^{(m)}, \mathbf{Z}_k} \quad (108a)$$

with, $P_{\Theta_k | \mathbf{U}_k^{(m)}, \mathbf{V}_k^{(m)}, \mathbf{Z}_k}$ defined in (28).

Furthermore, from Bayes' theorem [66], the joint probability measure $P_{\underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}$ in (107c) can be rewritten as,

$$P_{\underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0} = \prod_{t=1}^m P_{\mathbf{U}^{(t)}, \mathbf{V}^{(t)} | \underline{\mathbf{U}}^{(t-1)}, \underline{\mathbf{V}}^{(t-1)}, \mathbf{Z}_0} \quad (109)$$

$$= \prod_{t=1}^m P_{\mathbf{U}^{(t)} | \underline{\mathbf{U}}^{(t-1)}, \underline{\mathbf{V}}^{(t-1)}, \mathbf{Z}_0} P_{\mathbf{V}^{(t)} | \underline{\mathbf{U}}^{(t-1)}, \underline{\mathbf{V}}^{(t-1)}, \mathbf{Z}_0} \quad (110)$$

$$= \prod_{t=1}^m \left(P_{\mathbf{U}^{(t)} | \underline{\mathbf{U}}^{(t-1)}, \underline{\mathbf{V}}^{(t-1)}, \mathbf{Z}_0} \prod_{k=1}^K P_{\mathbf{V}_{k,t} | \mathbf{U}_k^{(t-1)}, \mathbf{V}_k^{(t-1)}, \mathbf{Z}_k} \right). \quad (111)$$

In the case in which $t = 1$, the notations $P_{\mathbf{U}^{(1)}, \mathbf{V}^{(1)} | \underline{\mathbf{U}}^{(0)}, \underline{\mathbf{V}}^{(0)}, \mathbf{Z}_0}$, $P_{\mathbf{U}^{(1)} | \underline{\mathbf{U}}^{(0)}, \underline{\mathbf{V}}^{(0)}, \mathbf{Z}_0}$, $P_{\mathbf{V}^{(1)} | \underline{\mathbf{U}}^{(0)}, \underline{\mathbf{V}}^{(0)}, \mathbf{Z}_0}$, $P_{\mathbf{U}^{(1)} | \underline{\mathbf{U}}^{(0)}, \underline{\mathbf{V}}^{(1)}, \mathbf{Z}_0}$, and $P_{\mathbf{V}_{k,1} | \mathbf{U}_k^{(0)}, \mathbf{V}_k^{(0)}, \mathbf{Z}_k}$ are longhands for the measures

$$P_{\mathbf{U}^{(1)}, \mathbf{V}^{(1)} | \mathbf{Z}_0} \in \Delta(\mathcal{U}_0 \times \mathcal{V}_0 | \mathcal{Z}_0) \quad (112a)$$

$$P_{\mathbf{U}^{(1)} | \underline{\mathbf{V}}^{(1)}, \mathbf{Z}_0} \in \Delta(\mathcal{U}_0 | \mathcal{V}_0 \times \mathcal{Z}_0), \quad (112b)$$

$$P_{\mathbf{V}^{(1)} | \mathbf{Z}_0} \in \Delta(\mathcal{V}_0 | \mathcal{Z}_0), \quad (112c)$$

$$P_{\mathbf{U}^{(1)} | \underline{\mathbf{V}}^{(1)}} \in \Delta(\mathcal{U}_0 | \mathcal{V}_0) \text{ and} \quad (112d)$$

$$P_{\mathbf{V}_{k,1} | \mathbf{Z}_k} \in \Delta(\mathcal{V}_k | \mathcal{Z}_k^{n_k}), \quad (112e)$$

respectively, and for all $t \in \{2, 3, \dots, m\}$,

$$P_{\mathbf{U}^{(t)}, \mathbf{V}^{(t)} | \underline{\mathbf{U}}^{(t-1)}, \underline{\mathbf{V}}^{(t-1)}, \mathbf{Z}_0} \in \Delta(\mathcal{U}_0 \times \mathcal{V}_0 | \mathcal{U}_0^{t-1} \times \mathcal{V}_0^{t-1} \times \mathcal{Z}_0), \quad (112f)$$

$$P_{\mathbf{U}^{(t)}|\underline{\mathbf{U}}^{(t-1)},\underline{\mathbf{V}}^{(t)},\mathbf{Z}_0} \in \Delta \left(\mathcal{U}_0 | \mathcal{U}_0^{t-1} \times \mathcal{V}_0^t \times \mathcal{Z}_0 \right), \quad (112g)$$

$$P_{\mathbf{V}^{(t)}|\underline{\mathbf{U}}^{(t-1)},\underline{\mathbf{V}}^{(t-1)},\mathbf{Z}_0} \in \Delta \left(\mathcal{V}_0 | \mathcal{U}_0^{t-1} \times \mathcal{V}_0^{t-1} \times \mathcal{Z}_0 \right). \quad (112h)$$

The conditional probability measure $P_{\mathbf{U}^{(t)}|\underline{\mathbf{U}}^{(t-1)},\underline{\mathbf{V}}^{(t)},\mathbf{Z}_0}$ in (110) simplifies to $P_{\mathbf{U}^{(t)}|\underline{\mathbf{U}}^{(t-1)},\underline{\mathbf{V}}^{(t)}}$ in (30) under Assumption 3. This is because the server does not have access to the training datasets \mathbf{z}_0 in (2). Similarly, under Assumption 2, the joint conditional probability $P_{\mathbf{V}^{(t)}|\underline{\mathbf{U}}^{(t-1)},\underline{\mathbf{V}}^{(t-1)},\mathbf{Z}_0}$ in (112h) satisfies

$$P_{\mathbf{V}^{(t)}|\underline{\mathbf{U}}^{(t-1)},\underline{\mathbf{V}}^{(t-1)},\mathbf{Z}_0} = \prod_{k=1}^K P_{V_{k,t}|\mathbf{U}_k^{(t-1)},\mathbf{V}_k^{(t-1)},\mathbf{Z}_k}, \quad (113)$$

with the conditional probability measure $P_{V_{k,t}|\mathbf{U}_k^{(t-1)},\mathbf{V}_k^{(t-1)},\mathbf{Z}_k}$ defined in (29). This completes the proof.

B Proof of Lemma 2

Consider the measurable set $\mathcal{A} \subseteq \mathcal{M}_k$, for some $k \in \{1, 2, \dots, K\}$, and define the set $\mathcal{A} \times \mathcal{M}_{-k} \subseteq \mathcal{M}_0$ as follows:

$$\mathcal{A} \times \mathcal{M}_{-k} \triangleq \mathcal{M}_1 \times \dots \times \mathcal{M}_{k-1} \times \mathcal{A} \times \mathcal{M}_{k+1} \times \dots \times \mathcal{M}_K. \quad (114)$$

Then, for all $\underline{\mathbf{u}}^{(m)} \in \mathcal{U}_0^m$, for all $\underline{\mathbf{v}}^{(m)} \in \mathcal{V}_0^m$ and for all $\mathbf{z}_0 \in \mathcal{Z}_0$

$$P_{\Theta|\underline{\mathbf{U}}^{(m)}=\underline{\mathbf{u}}^{(m)},\underline{\mathbf{V}}^{(m)}=\underline{\mathbf{v}}^{(m)},\mathbf{Z}_0=\mathbf{z}_0}(\mathcal{A} \times \mathcal{M}_{-k}) = \prod_{t=1}^{k-1} \left(P_{\Theta_t|\mathbf{U}_t^{(m)}=\mathbf{u}_t^{(m)},\mathbf{V}_t^{(m)}=\mathbf{v}_t^{(m)},\mathbf{Z}_t=\mathbf{z}_t}(\mathcal{M}_t) \right)$$

$$P_{\Theta_k|\mathbf{U}_k^{(m)}=\mathbf{u}_k^{(m)},\mathbf{V}_k^{(m)}=\mathbf{v}_k^{(m)},\mathbf{Z}_k=\mathbf{z}_k}(\mathcal{A}) \prod_{t=k+1}^K \left(P_{\Theta_t|\mathbf{U}_t^{(m)}=\mathbf{u}_t^{(m)},\mathbf{V}_t^{(m)}=\mathbf{v}_t^{(m)},\mathbf{Z}_t=\mathbf{z}_t}(\mathcal{M}_t) \right) \quad (115)$$

$$= P_{\Theta_k|\mathbf{U}_k^{(m)}=\mathbf{u}_k^{(m)},\mathbf{V}_k^{(m)}=\mathbf{v}_k^{(m)},\mathbf{Z}_k=\mathbf{z}_k}(\mathcal{A}), \quad (116)$$

where the equality in (115) follows from Assumption 1. Moreover, for all $i \in \{1, \dots, K\}$, the conditional measure $P_{\Theta_i|\mathbf{U}_i^{(m)},\mathbf{V}_i^{(m)},\mathbf{Z}_i}$ in (31) satisfies, for all $\mathbf{u}_i^{(m)} \in \mathcal{U}_i^m$, for all $\mathbf{v}_i^{(m)} \in \mathcal{V}_i^m$ and for all $\mathbf{z}_i \in \mathcal{Z}_i^{n_i}$,

$$P_{\Theta_i|\mathbf{U}_i^{(m)}=\mathbf{u}_i^{(m)},\mathbf{V}_i^{(m)}=\mathbf{v}_i^{(m)},\mathbf{Z}_i=\mathbf{z}_i}(\mathcal{M}_i) = 1, \quad (117)$$

which justifies (116).

On the other hand, for all $\mathbf{z}_0 \in \mathcal{Z}_0$, it follows that

$$\begin{aligned} & \int P_{\Theta|\underline{\mathbf{U}}^{(m)}=\underline{\mathbf{u}}^{(m)},\underline{\mathbf{V}}^{(m)}=\underline{\mathbf{v}}^{(m)},\mathbf{Z}_0=\mathbf{z}_0}(\mathcal{A} \times \mathcal{M}_{-k}) dP_{\underline{\mathbf{U}}^{(m)},\underline{\mathbf{V}}^{(m)}|\mathbf{Z}_0=\mathbf{z}_0}(\underline{\mathbf{u}}^{(m)},\underline{\mathbf{v}}^{(m)}) \\ &= \int P_{\Theta_k|\mathbf{U}_k^{(m)}=\mathbf{u}_k^{(m)},\mathbf{V}_k^{(m)}=\mathbf{v}_k^{(m)},\mathbf{Z}_k=\mathbf{z}_k}(\mathcal{A}) dP_{\underline{\mathbf{U}}^{(m)},\underline{\mathbf{V}}^{(m)}|\mathbf{Z}_0=\mathbf{z}_0}(\underline{\mathbf{u}}^{(m)},\underline{\mathbf{v}}^{(m)}) \end{aligned} \quad (118)$$

$$= P_{\Theta_k | Z_0 = z_0}^{(m)}(\mathcal{A}), \quad (119)$$

where the equality in (118) follows from (116); and the equality in (119) represents the marginal probability measure in $\Delta(\mathcal{M}_k | Z_0)$ of the probability measure $P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}$ in (26).

Finally, given that the sets \mathcal{U}_0 and \mathcal{V}_0 are finite, the integral in (118) reduces to a finite sum,

$$\begin{aligned} & P_{\Theta_k | Z_0 = z_0}^{(m)}(\mathcal{A}) \\ &= \int P_{\Theta_k | U_k^{(m)} = \underline{u}_k^{(m)}, V_k^{(m)} = \underline{v}_k^{(m)}, Z_k = z_k}(\mathcal{A}) dP_{\underline{U}^{(m)}, \underline{V}^{(m)} | Z_0 = z_0}(\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}) \end{aligned} \quad (120)$$

$$= \sum_{(\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}) \in \mathcal{U}_0^m \times \mathcal{V}_0^m} a_{\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, z_0} P_{\Theta_k | U_k^{(m)} = \underline{u}_k^{(m)}, V_k^{(m)} = \underline{v}_k^{(m)}, Z_k = z_k}(\mathcal{A}), \quad (121)$$

where the coefficients $a_{\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, z_0}$ satisfy

$$a_{\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, z_0} \triangleq P_{\underline{U}^{(m)}, \underline{V}^{(m)} | Z_0 = z_0}(\{(\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)})\}) \quad (122)$$

$$\begin{aligned} &= \prod_{t=1}^m P_{\underline{U}^{(t)} | \underline{U}^{(t-1)} = \underline{\mathbf{u}}^{(t-1)}, \underline{V}^{(t)} = \underline{\mathbf{v}}^{(t)}}(\{\underline{\mathbf{u}}^{(t)}\}) \\ & \quad \prod_{j=1}^K P_{V_{j,t} | U_j^{(t-1)} = \underline{u}_j^{(t-1)}, V_j^{(t-1)} = \underline{v}_j^{(t-1)}, Z_j = z_j}(\{v_{j,t}\}), \end{aligned} \quad (123)$$

where (123) follows from (111). Note that from (122), it follows that $a_{\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, z_0} \geq 0$; and

$$\sum_{(\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}) \in \mathcal{U}_0^m \times \mathcal{V}_0^m} a_{\underline{\mathbf{u}}^{(m)}, \underline{\mathbf{v}}^{(m)}, z_0} = 1, \quad (124)$$

which completes the proof.

C Proof of Lemma 3

The proof of (57) is primarily algebraic and starts as follows

$$\begin{aligned} & \int R_{k, z_k} \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \right) dP_{Z_0}(z_0) \\ &= \iint L_k(z_k, \theta_k) dP_{\Theta_k | Z_0 = z_0}^{(m)}(\theta_k) dP_{Z_0}(z_0) \end{aligned} \quad (125)$$

$$= \iint L_k(z_k, \theta_k) \frac{dP_{\Theta_k | Z_0 = z_0}^{(m)}(\theta_k)}{dP_{\Theta_k}^{(m)}(\theta_k)} dP_{\Theta_k}^{(m)}(\theta_k) dP_{Z_0}(z_0) \quad (126)$$

$$= \iint L_k(z_k, \theta_k) \frac{dP_{Z_0 | \Theta_k = \theta_k}^{(m)}(z_0)}{dP_{Z_0}^{(m)}(z_0)} dP_{\Theta_k}^{(m)}(\theta_k) dP_{Z_0}(z_0) \quad (127)$$

$$= \iint \mathsf{L}_k(\mathbf{z}_k, \boldsymbol{\theta}_k) \frac{dP_{\mathbf{Z}_0|\boldsymbol{\Theta}_k=\boldsymbol{\theta}_k}^{(m)}(\mathbf{z}_0)}{dP_{\mathbf{Z}_0}}(\mathbf{z}_0) dP_{\mathbf{Z}_0}(\mathbf{z}_0) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) \quad (128)$$

$$= \iint \mathsf{L}_k(\mathbf{z}_k, \boldsymbol{\theta}_k) dP_{\mathbf{Z}_0|\boldsymbol{\Theta}_k=\boldsymbol{\theta}_k}^{(m)}(\mathbf{z}_0) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) \quad (129)$$

$$= \iint \mathsf{L}_k(\mathbf{z}_k, \boldsymbol{\theta}_k) dP_{\mathbf{Z}_k|\boldsymbol{\Theta}_k=\boldsymbol{\theta}_k}^{(m)}(\mathbf{z}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) \quad (130)$$

$$= \int \mathsf{R}_{k,\boldsymbol{\theta}_k} \left(P_{\mathbf{Z}_k|\boldsymbol{\Theta}_k=\boldsymbol{\theta}_k}^{(m)} \right) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k), \quad (131)$$

where (126) follows from [58, Theorem 2]; (127) follows from [58, Theorem 11]; (128) follows from follows from [67, Theorem 2.6.6]; (129) follows from [58, Theorem 2]; (130) from the fact that the functional L_k depends exclusively on \mathbf{z}_k (instead of \mathbf{z}_0); and (131) follows from the definition of the functional $\mathsf{R}_{k,\boldsymbol{\theta}_k}$ in (53). This completes the proof.

D Proof of Lemma 4

The proof of Lemma 4 is primarily algebraic and proceeds as follows

$$\begin{aligned} & \iint \mathsf{R}_{k,\widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)} \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) \\ &= \iiint \mathsf{L}_k(\widehat{\mathbf{z}}_k, \boldsymbol{\theta}_k) dP_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}(\boldsymbol{\theta}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) \end{aligned} \quad (132)$$

$$= \iiint \mathsf{L}_k(\widehat{\mathbf{z}}_k, \boldsymbol{\theta}_k) \frac{dP_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}(\boldsymbol{\theta}_k)}{dP_{\boldsymbol{\Theta}_k}^{(m)}}(\boldsymbol{\theta}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) \quad (133)$$

$$= \iint \mathsf{L}_k(\widehat{\mathbf{z}}_k, \boldsymbol{\theta}_k) \int \frac{dP_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}(\boldsymbol{\theta}_k)}{dP_{\boldsymbol{\Theta}_k}^{(m)}}(\boldsymbol{\theta}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) \quad (134)$$

$$= \iint \mathsf{L}_k(\widehat{\mathbf{z}}_k, \boldsymbol{\theta}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) \quad (135)$$

$$= \int \mathsf{R}_{k,\widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\Theta}_k}^{(m)} \right) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) \quad (136)$$

$$= \int \mathsf{R}_{k,\boldsymbol{\theta}_k} \left(P_{\mathbf{Z}_k} \right) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k), \quad (137)$$

where (133) follows from a change of measure [58, Theorem 2]; (134) follows from [67, Theorem 2.6.6]; (135) follows from [58, Theorem 10]; and (136) follows from the definition of the functional $\mathsf{R}_{k,\widehat{\mathbf{z}}_k}$ in (52) and concludes the proof of (58). The equality in (137) follows from [67, Theorem 2.6.6] and the definition of the functional $\mathsf{R}_{k,\boldsymbol{\theta}_k}$ in (53). This concludes the proof of (59). The proof continues from (136) as follows

$$\int \mathsf{R}_{k,\widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\Theta}_k}^{(m)} \right) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k)$$

$$= \iint \mathbb{L}_k(\widehat{\mathbf{z}}_k, \boldsymbol{\theta}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) \quad (138)$$

$$= \iint \mathbb{L}_k(\mathbf{z}_k, \widehat{\boldsymbol{\theta}}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\widehat{\boldsymbol{\theta}}_k) \int \frac{dP_{\mathbf{Z}_k|\boldsymbol{\Theta}_k=\boldsymbol{\theta}_k}^{(m)}(\mathbf{z}_k)}{dP_{\mathbf{Z}_k}}(\mathbf{z}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) dP_{\mathbf{Z}_k}(\mathbf{z}_k) \quad (139)$$

$$= \iiint \mathbb{L}_k(\mathbf{z}_k, \widehat{\boldsymbol{\theta}}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\widehat{\boldsymbol{\theta}}_k) \frac{dP_{\mathbf{Z}_k|\boldsymbol{\Theta}_k=\boldsymbol{\theta}_k}^{(m)}(\mathbf{z}_k)}{dP_{\mathbf{Z}_k}}(\mathbf{z}_k) dP_{\mathbf{Z}_k}(\mathbf{z}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) \quad (140)$$

$$= \iiint \mathbb{L}_k(\mathbf{z}_k, \widehat{\boldsymbol{\theta}}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\widehat{\boldsymbol{\theta}}_k) dP_{\mathbf{Z}_k|\boldsymbol{\Theta}_k=\boldsymbol{\theta}_k}^{(m)}(\mathbf{z}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) \quad (141)$$

$$= \iint \mathbb{R}_{k,\widehat{\boldsymbol{\theta}}_k}(P_{\mathbf{Z}_k|\boldsymbol{\Theta}_k=\boldsymbol{\theta}_k}^{(m)}) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\widehat{\boldsymbol{\theta}}_k), \quad (142)$$

where (139) follows from [58, Theorem 10]; (140) from [67, Theorem 2.6.6]; and (141) from [58, Theorem 2]. The equality in (142) concludes the proof of (60).

E Proof of Lemma 5

The proof follows from Definition 12, from which it holds that

$$\begin{aligned} & \mathbf{G}_k(P_{\underline{\boldsymbol{\Theta}}, \underline{\mathbf{V}}^{(m)}, \underline{\mathbf{V}}^{(m)}|\mathbf{Z}_0}; P_{\mathbf{Z}_0}) \\ &= \iint \left(\mathbb{R}_{k,\widehat{\mathbf{z}}_k}(P_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}) - \mathbb{R}_{k,\mathbf{z}_k}(P_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}) \right) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) \end{aligned} \quad (143)$$

$$= \iint \mathbb{R}_{k,\widehat{\mathbf{z}}_k}(P_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) - \int \mathbb{R}_{k,\mathbf{z}_k}(P_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}) dP_{\mathbf{Z}_0}(\mathbf{z}_0). \quad (144)$$

The proof focuses on the term $\iint \mathbb{R}_{k,\widehat{\mathbf{z}}_k}(P_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0)$ in (144). From Lemma 4, it follows that

$$\begin{aligned} & \iint \mathbb{R}_{k,\widehat{\mathbf{z}}_k}(P_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) \\ &= \iint \mathbb{R}_{k,\widehat{\boldsymbol{\theta}}_k}(P_{\mathbf{Z}_k|\boldsymbol{\Theta}_k=\boldsymbol{\theta}_k}^{(m)}) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k) dP_{\boldsymbol{\Theta}_k}^{(m)}(\widehat{\boldsymbol{\theta}}_k). \end{aligned} \quad (145)$$

Consider now the term $\int \mathbb{R}_{k,\mathbf{z}_k}(P_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}) dP_{\mathbf{Z}_0}(\mathbf{z}_0)$ in (144). From Lemma 3, it follows that

$$\int \mathbb{R}_{k,\mathbf{z}_k}(P_{\boldsymbol{\Theta}_k|\mathbf{Z}_0=\mathbf{z}_0}^{(m)}) dP_{\mathbf{Z}_0}(\mathbf{z}_0) = \int \mathbb{R}_{k,\boldsymbol{\theta}_k}(P_{\mathbf{Z}_k|\boldsymbol{\Theta}_k=\boldsymbol{\theta}_k}^{(m)}) dP_{\boldsymbol{\Theta}_k}^{(m)}(\boldsymbol{\theta}_k). \quad (146)$$

The proof is completed by using (145) and (146) in (144).

F Proof of Lemma 12

The proof of (86) is primarily algebraic and proceeds by analyzing the KL divergence $D(P_1 \parallel P_{\boldsymbol{\Theta}_k|\mathbf{Z}_k=\mathbf{z}_k}^{(\mathcal{Q}_{\boldsymbol{\Theta}_k}, \lambda_k)})$. More precisely,

$$D(P_1 \parallel P_{\boldsymbol{\Theta}_k|\mathbf{Z}_k=\mathbf{z}_k}^{(\mathcal{Q}_{\boldsymbol{\Theta}_k}, \lambda_k)})$$

$$= \int \log \left(\frac{dP_1}{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}}(\boldsymbol{\theta}) \right) dP_1(\boldsymbol{\theta}) \quad (147)$$

$$= \int \log \left(\frac{dQ_{\Theta_k}}{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}}(\boldsymbol{\theta}_k) \frac{dP_1}{dQ_{\Theta_k}}(\boldsymbol{\theta}_k) \right) dP_1(\boldsymbol{\theta}_k) \quad (148)$$

$$= \int \log \left(\left(\frac{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)}}{dQ_{\Theta_k}}(\boldsymbol{\theta}_k) \right)^{-1} \frac{dP_1}{dQ_{\Theta_k}}(\boldsymbol{\theta}_k) \right) dP_1(\boldsymbol{\theta}_k) \quad (149)$$

$$= \int -\log \left(\frac{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)}}{dQ_{\Theta_k}}(\boldsymbol{\theta}_k) \right) + \log \left(\frac{dP_1}{dQ_{\Theta_k}}(\boldsymbol{\theta}_k) \right) dP_1(\boldsymbol{\theta}_k) \quad (150)$$

$$= - \int \log \left(\frac{dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)}}{dQ_{\Theta_k}}(\boldsymbol{\theta}_k) \right) dP_1(\boldsymbol{\theta}_k) + D(P_1 \parallel Q_{\Theta_k}) \quad (151)$$

$$= - \int \left(-\frac{1}{\lambda_k} \mathsf{L}_k(\mathbf{z}_k, \boldsymbol{\theta}_k) - \mathsf{K}_{k, Q_{\Theta_k}, \mathbf{z}_k} \left(-\frac{1}{\lambda_k} \right) \right) dP_1(\boldsymbol{\theta}_k) + D(P_1 \parallel Q_{\Theta_k}) \quad (152)$$

$$= \int \frac{1}{\lambda_k} \mathsf{L}_k(\mathbf{z}_k, \boldsymbol{\theta}_k) dP_1(\boldsymbol{\theta}_k) + \mathsf{K}_{k, Q_{\Theta_k}, \mathbf{z}_k} \left(-\frac{1}{\lambda_k} \right) + D(P_1 \parallel Q_{\Theta_k}) \quad (153)$$

$$= \int \frac{1}{\lambda_k} \mathsf{L}_k(\mathbf{z}_k, \boldsymbol{\theta}_k) dP_1(\boldsymbol{\theta}_k) - \frac{1}{\lambda_k} \int \mathsf{L}_k(\mathbf{z}_k, \boldsymbol{\theta}_k) dP_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)}(\boldsymbol{\theta}_k) - D \left(P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \parallel Q_{\Theta_k} \right) + D(P_1 \parallel Q_{\Theta_k}) \quad (154)$$

$$= \frac{1}{\lambda_k} \left(\mathsf{R}_{k, \mathbf{z}_k}(P_1) - \mathsf{R}_{k, \mathbf{z}_k} \left(P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right) - D \left(P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \parallel Q_{\Theta_k} \right) + D(P_1 \parallel Q_{\Theta_k}) \quad (155)$$

$$= \frac{1}{\lambda_k} \mathsf{G}_k(\mathbf{z}_k, P_1, P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)}) - D \left(P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \parallel Q_{\Theta_k} \right) + D(P_1 \parallel Q_{\Theta_k}), \quad (156)$$

where equality (148) follows from [58, Theorem 4]; (149) from [58, Theorem 5]; (152) from (66); (154) from Lemma 7; and (156) from (82).

Hence, from the equality in (156), it follows that

$$\begin{aligned} & \mathsf{G}_k \left(\mathbf{z}_k, P_1, P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \\ &= \lambda_k \left(D \left(P_1 \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) + D \left(P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \parallel Q_{\Theta_k} \right) - D(P_1 \parallel Q_{\Theta_k}) \right). \end{aligned} \quad (157)$$

The proof proceeds from the definition of the functional G_k in (82), from which it holds that

$$\begin{aligned} & \mathsf{G}_k(\mathbf{z}_k, P_1, P_2) \\ &= \mathsf{R}_{k, \mathbf{z}_k}(P_1) - \mathsf{R}_{k, \mathbf{z}_k}(P_2) + \mathsf{R}_{k, \mathbf{z}_k} \left(P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) - \mathsf{R}_{k, \mathbf{z}_k} \left(P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \end{aligned} \quad (158)$$

$$= \mathsf{G}_k \left(\mathbf{z}_k, P_1, P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) - \mathsf{G}_k \left(\mathbf{z}_k, P_2, P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \quad (159)$$

$$= \lambda_k \left(D \left(P_1 \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) + D \left(P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \parallel Q_{\Theta_k} \right) - D(P_1 \parallel Q_{\Theta_k}) \right)$$

$$-\lambda_k \left(D \left(P_2 \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) + D \left(P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \parallel Q_{\Theta_k} \right) - D \left(P_2 \parallel Q_{\Theta_k} \right) \right) \quad (160)$$

$$= \lambda_k \left(D \left(P_1 \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) - D \left(P_2 \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) + D \left(P_2 \parallel Q_{\Theta_k} \right) - D \left(P_1 \parallel Q_{\Theta_k} \right) \right), \quad (161)$$

where the equality in (159) follows from the definition of the functional G_k in (82); and the equality in (160) follows from (157). This completes the proof.

G Proof of Lemma 13

The proof follows along the same lines as the proof of Lemma 12. It is included entirely hereunder for the sake of completeness.

The proof of (87) is primarily algebraic and proceeds by analyzing the KL divergence $D \left(P_1 \parallel P_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right)$. More precisely,

$$\begin{aligned} & D \left(P_1 \parallel P_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \\ &= \int \log \left(\frac{dP_1}{dP_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}} (\mathbf{z}_k) \right) dP_1 (\mathbf{z}_k) \end{aligned} \quad (162)$$

$$= \int \log \left(\frac{dQ_{\mathbf{Z}_k}}{dP_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}} (\mathbf{z}_k) \frac{dP_1}{dQ_{\mathbf{Z}_k}} (\mathbf{z}_k) \right) dP_1 (\mathbf{z}_k) \quad (163)$$

$$= \int \log \left(\left(\frac{dP_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}{dQ_{\mathbf{Z}_k}} (\mathbf{z}_k) \right)^{-1} \frac{dP_1}{dQ_{\mathbf{Z}_k}} (\mathbf{z}_k) \right) dP_1 (\mathbf{z}_k) \quad (164)$$

$$= \int -\log \left(\frac{dP_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}{dQ_{\mathbf{Z}_k}} (\mathbf{z}_k) \right) + \log \left(\frac{dP_1}{dQ_{\mathbf{Z}_k}} (\mathbf{z}_k) \right) dP_1 (\mathbf{z}_k) \quad (165)$$

$$= - \int \log \left(\frac{dP_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}{dQ_{\mathbf{Z}_k}} (\mathbf{z}_k) \right) dP_1 (\mathbf{z}_k) + D \left(P_1 \parallel Q_{\mathbf{Z}_k} \right) \quad (166)$$

$$= - \int \left(-\frac{1}{\alpha_k} \mathsf{L}_k (\mathbf{z}_k, \theta_k) - \mathsf{K}_{k, Q_{\mathbf{Z}_k}, \theta_k} \left(-\frac{1}{\alpha_k} \right) \right) dP_1 (\mathbf{z}_k) + D \left(P_1 \parallel Q_{\mathbf{Z}_k} \right) \quad (167)$$

$$= \mathsf{K}_{k, Q_{\mathbf{Z}_k}, \theta_k} \left(-\frac{1}{\alpha_k} \right) + \int \frac{1}{\alpha_k} \mathsf{L}_k (\mathbf{z}_k, \theta_k) dP_1 (\mathbf{z}_k) + D \left(P_1 \parallel Q_{\mathbf{Z}_k} \right) \quad (168)$$

$$\begin{aligned} &= -\frac{1}{\alpha_k} \int \mathsf{L}_k (\mathbf{z}_k, \theta_k) dP_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} (\mathbf{z}_k) - D \left(P_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) \\ &\quad + \int \frac{1}{\alpha_k} \mathsf{L}_k (\mathbf{z}_k, \theta_k) dP_1 (\mathbf{z}_k) + D \left(P_1 \parallel Q_{\mathbf{Z}_k} \right) \end{aligned} \quad (169)$$

$$= \frac{1}{\alpha_k} \left(\mathsf{R}_{k, \theta_k} (P_1) - \mathsf{R}_{k, \theta_k} \left(P_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \right) - D \left(P_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) + D \left(P_1 \parallel Q_{\mathbf{Z}_k} \right) \quad (170)$$

$$= \frac{1}{\alpha_k} \mathbf{G}_k \left(\boldsymbol{\theta}_k, P_1, P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) - D \left(P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) + D(P_1 \parallel Q_{\mathbf{Z}_k}), \quad (171)$$

where equality (163) follows from [58, Theorem 4]; (164) from [58, Theorem 5]; (167) from (66); (169) from Lemma 9; and (171) from (82).

Hence, from the equality in (171) it follows that

$$\begin{aligned} & \mathbf{G}_k \left(\boldsymbol{\theta}_k, P_1, P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \\ &= \alpha_k \left(D \left(P_1 \parallel P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) + D \left(P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) - D(P_1 \parallel Q_{\mathbf{Z}_k}) \right). \end{aligned} \quad (172)$$

The proof proceeds from the definition of the functional \mathbf{G}_k in (83), from which it holds that

$$\begin{aligned} & \mathbf{G}_k(\boldsymbol{\theta}_k, P_1, P_2) \\ &= \mathbf{R}_{k, \boldsymbol{\theta}_k}(P_1) - \mathbf{R}_{k, \boldsymbol{\theta}_k}(P_2) + \mathbf{R}_{k, \boldsymbol{\theta}_k} \left(P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) - \mathbf{R}_{k, \boldsymbol{\theta}_k} \left(P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \end{aligned} \quad (173)$$

$$= \mathbf{G}_k \left(\boldsymbol{\theta}_k, P_1, P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) - \mathbf{G}_k \left(\boldsymbol{\theta}_k, P_2, P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \quad (174)$$

$$\begin{aligned} &= \alpha_k \left(D \left(P_1 \parallel P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) + D \left(P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) - D(P_1 \parallel Q_{\mathbf{Z}_k}) \right) \\ &\quad - \alpha_k \left(D \left(P_2 \parallel P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) + D \left(P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) - D(P_2 \parallel Q_{\mathbf{Z}_k}) \right) \end{aligned} \quad (175)$$

$$= \alpha_k \left(D \left(P_1 \parallel P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) - D(P_1 \parallel Q_{\mathbf{Z}_k}) - D \left(P_2 \parallel P_{\widehat{\mathbf{Z}}_k | \boldsymbol{\Theta}_k = \boldsymbol{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) + D(P_2 \parallel Q_{\mathbf{Z}_k}) \right), \quad (176)$$

where the equality in (174) follows from the definition of the functional \mathbf{G}_k in (83); and the equality in (175) follows from (172). This completes the proof.

H Proof of Lemma 14

The proof follows from Definition 12, which yields

$$\begin{aligned} & \mathbf{G}_k \left(P_{\underline{\boldsymbol{\theta}}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) \\ &= \iint \left(\mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) - \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) \right) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) \end{aligned} \quad (177)$$

$$= \iint \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) - \int \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0). \quad (178)$$

The proof focuses on the term $\iint \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0)$ in (144). From Lemma 4, it follows that

$$\iint \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\Theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) = \int \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\Theta}_k}^{(m)} \right) dP_{\mathbf{Z}_k}(\mathbf{z}_k). \quad (179)$$

Using (179) in (178) yields

$$\mathbf{G}_k \left(P_{\underline{\boldsymbol{\theta}}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$$

$$= \int \mathbb{R}_{k,z_k} \left(P_{\Theta_k}^{(m)} \right) dP_{Z_k} (z_k) - \int \mathbb{R}_{k,z_k} \left(P_{\Theta_k|Z_0=z_0}^{(m)} \right) dP_{Z_0} (z_0) \quad (180)$$

$$= \int \mathbb{R}_{k,z_k} \left(P_{\Theta_k}^{(m)} \right) - \mathbb{R}_{k,z_k} \left(P_{\Theta_k|Z_0=z_0}^{(m)} \right) dP_{Z_0} (z_0) \quad (181)$$

$$= \int \mathbb{G}_k \left(z_k; P_{\Theta_k}^{(m)}; P_{\Theta_k|Z_0=z_0}^{(m)} \right) dP_{Z_0} (z_0), \quad (182)$$

which completes the proof of (88).

I Proof of Lemma 15

The proof follows along the same lines as the proof of Lemma 14. It is included entirely hereunder for the sake of completeness.

The proof follows from Lemma 5, which yields

$$\begin{aligned} & \mathbb{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0}; P_{Z_0} \right) \\ &= \iint \left(\mathbb{R}_{k,\hat{\theta}_k} \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \right) - \mathbb{R}_{k,\theta_k} \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \right) \right) dP_{\Theta_k}^{(m)} \left(\hat{\theta}_k \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right) \end{aligned} \quad (183)$$

$$= \iint \mathbb{R}_{k,\hat{\theta}_k} \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \right) dP_{\Theta_k}^{(m)} \left(\hat{\theta}_k \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right) - \int \mathbb{R}_{k,\theta_k} \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right). \quad (184)$$

The proof focuses on the term $\iint \mathbb{R}_{k,\hat{\theta}_k} \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \right) dP_{\Theta_k}^{(m)} \left(\hat{\theta}_k \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right)$ in (184). From Lemma 4, it follows that

$$\iint \mathbb{R}_{k,\hat{\theta}_k} \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \right) dP_{\Theta_k}^{(m)} \left(\hat{\theta}_k \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right) = \int \mathbb{R}_{k,\theta_k} \left(P_{Z_k} \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right). \quad (185)$$

Using (185) in (184) yields

$$\begin{aligned} & \mathbb{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0}; P_{Z_0} \right) \\ &= \int \mathbb{R}_{k,\theta_k} \left(P_{Z_k} \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right) - \int \mathbb{R}_{k,\theta_k} \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right) \end{aligned} \quad (186)$$

$$= \int \mathbb{R}_{k,\theta_k} \left(P_{Z_k} \right) - \mathbb{R}_{k,\theta_k} \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right) \quad (187)$$

$$= \int \mathbb{G}_k \left(\theta_k, P_{Z_k}, P_{Z_k|\Theta_k=\theta_k}^{(m)} \right) dP_{\Theta_k}^{(m)} \left(\theta_k \right), \quad (188)$$

which completes the proof of (89).

J Proof of Theorem 18

The proof starts by considering the information measure $J_{L_k, Q_{\theta_k}, \lambda_k} \left(P_{\theta_k|Z_0}^{(m)}; P_{Z_0} \right)$ in (92), from which it follows that

$$\begin{aligned}
J_{L_k, Q_{\theta_k}, \lambda_k} \left(P_{\theta_k|Z_0}^{(m)}; P_{Z_0} \right) &= \int_{\mathcal{A}_\gamma} \log \left(\frac{dP_{\theta_k|Z_0=z_0}^{(m)}}{dP_{\theta_k|Z_k=z_k}^{(Q_{\theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\theta_k}^{(m)} P_{Z_0} (\theta_k, z_0) \\
&+ \int_{\mathcal{A}_\gamma^c} \log \left(\frac{dP_{\theta_k|Z_0=z_0}^{(m)}}{dP_{\theta_k|Z_k=z_k}^{(Q_{\theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\theta_k}^{(m)} P_{Z_0} (\theta_k, z_0) \\
&- \int_{\mathcal{A}_\gamma} \log \left(\frac{dP_{\theta_k|Z_0=z_0}^{(m)}}{dP_{\theta_k|Z_k=z_k}^{(Q_{\theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\theta_k|Z_0}^{(m)} P_{Z_0} (\theta_k, z_0) \\
&- \int_{\mathcal{A}_\gamma^c} \log \left(\frac{dP_{\theta_k|Z_0=z_0}^{(m)}}{dP_{\theta_k|Z_k=z_k}^{(Q_{\theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\theta_k|Z_0}^{(m)} P_{Z_0} (\theta_k, z_0) \tag{189}
\end{aligned}$$

$$\begin{aligned}
&\leq \gamma P_{\theta_k}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma^c \right) - \gamma P_{\theta_k|Z_0}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma \right) \\
&+ \max_{(\theta_k, z_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{\theta_k|Z_0=z_0}^{(m)}}{dP_{\theta_k|Z_k=z_k}^{(Q_{\theta_k}, \lambda_k)}} (\theta_k) \right) P_{\theta_k}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma \right) \\
&- \min_{(\theta_k, z_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{\theta_k|Z_0=z_0}^{(m)}}{dP_{\theta_k|Z_k=z_k}^{(Q_{\theta_k}, \lambda_k)}} (\theta_k) \right) P_{\theta_k|Z_0}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma^c \right) \tag{190}
\end{aligned}$$

$$\begin{aligned}
&= -\gamma P_{\theta_k}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma \right) + \gamma P_{\theta_k|Z_0}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma^c \right) \\
&+ \max_{(\theta_k, z_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{\theta_k|Z_0=z_0}^{(m)}}{dP_{\theta_k|Z_k=z_k}^{(Q_{\theta_k}, \lambda_k)}} (\theta_k) \right) P_{\theta_k}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma \right) \\
&- \min_{(\theta_k, z_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{\theta_k|Z_0=z_0}^{(m)}}{dP_{\theta_k|Z_k=z_k}^{(Q_{\theta_k}, \lambda_k)}} (\theta_k) \right) P_{\theta_k|Z_0}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma^c \right) \tag{191}
\end{aligned}$$

$$\begin{aligned}
&= P_{\theta_k}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma \right) \left(\max_{(\theta_k, z_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{\theta_k|Z_0=z_0}^{(m)}}{dP_{\theta_k|Z_k=z_k}^{(Q_{\theta_k}, \lambda_k)}} (\theta_k) \right) - \gamma \right) \\
&+ P_{\theta_k|Z_0}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma^c \right) \left(\gamma - \min_{(\theta_k, z_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{\theta_k|Z_0=z_0}^{(m)}}{dP_{\theta_k|Z_k=z_k}^{(Q_{\theta_k}, \lambda_k)}} (\theta_k) \right) \right), \tag{192}
\end{aligned}$$

where the inequality in (190) follows from the definition of the set \mathcal{A}_γ and the extrema of the log-likelihood ratio.

A lower bounds on the information measure $J_{L_k, Q_{\Theta_k}, \lambda_k} \left(P_{\Theta_k|Z_0}^{(m)}; P_{Z_0} \right)$ in (92) can also be obtained as follows,

$$\begin{aligned}
J_{L_k, Q_{\Theta_k}, \lambda_k} \left(P_{\Theta_k|Z_0}^{(m)}; P_{Z_0} \right) &= \int_{\mathcal{A}_\gamma} \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\Theta_k}^{(m)} P_{Z_0} (\theta_k, z_0) \\
&+ \int_{\mathcal{A}_\gamma^c} \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\Theta_k}^{(m)} P_{Z_0} (\theta_k, z_0) \\
&- \int_{\mathcal{A}_\gamma} \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\Theta_k|Z_0}^{(m)} P_{Z_0} (\theta_k, z_0) \\
&- \int_{\mathcal{A}_\gamma^c} \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\Theta_k|Z_0}^{(m)} P_{Z_0} (\theta_k, z_0) \tag{193}
\end{aligned}$$

$$\begin{aligned}
&\geq \gamma P_{\Theta_k}^{(m)} P_{Z_0} (\mathcal{A}_\gamma) - \gamma P_{\Theta_k|Z_0}^{(m)} P_{Z_0} (\mathcal{A}_\gamma^c) \\
&+ \min_{(\theta_k, z_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) P_{\Theta_k}^{(m)} P_{Z_0} (\mathcal{A}_\gamma^c) \\
&- \max_{(\theta_k, z_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) P_{\Theta_k|Z_0}^{(m)} P_{Z_0} (\mathcal{A}_\gamma) \tag{194}
\end{aligned}$$

$$\begin{aligned}
&= -\gamma P_{\Theta_k}^{(m)} P_{Z_0} (\mathcal{A}_\gamma^c) + \gamma P_{\Theta_k|Z_0}^{(m)} P_{Z_0} (\mathcal{A}_\gamma) \\
&+ \min_{(\theta_k, z_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) P_{\Theta_k}^{(m)} P_{Z_0} (\mathcal{A}_\gamma^c) \\
&- \max_{(\theta_k, z_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) P_{\Theta_k|Z_0}^{(m)} P_{Z_0} (\mathcal{A}_\gamma) \tag{195}
\end{aligned}$$

$$\begin{aligned}
&= P_{\Theta_k}^{(m)} P_{Z_0} (\mathcal{A}_\gamma^c) \left(\min_{(\theta_k, z_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) - \gamma \right) \\
&+ P_{\Theta_k|Z_0}^{(m)} P_{Z_0} (\mathcal{A}_\gamma) \left(\gamma - \max_{(\theta_k, z_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) \right), \tag{196}
\end{aligned}$$

where the equality in (194) follows from the definition of the set \mathcal{A}_γ and the extrema of the log-likelihood ratio.

K Proof of Theorem 19

The following lemma is introduced as part of this proof.

Lemma 26. Consider the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$ in (62) and assume that for all $\mathbf{z}_0 \in \mathbf{Z}_0$, the probability measures $P_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}$ in (50); $P_{\Theta_k}^{(m)}$ in (55); and the σ -finite measure Q_{Θ_k} in (66) satisfy $P_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \ll Q_{\Theta_k} \ll P_{\Theta_k}^{(m)} \ll Q_{\Theta_k}$. Then,

$$\begin{aligned} \mathbf{G}_k \left(P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) &= \lambda_k I \left(P_{\Theta_k | \mathbf{Z}_0}^{(m)}; P_{\mathbf{Z}_0} \right) \\ &+ \lambda_k \int \left(D \left(P_{\Theta_k}^{(m)} \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) - D \left(P_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0), \end{aligned} \quad (197)$$

where the conditional probability measure $P_{\Theta_k | \mathbf{Z}_k}^{(Q_{\Theta_k}, \lambda_k)} \in \Delta \left(\mathcal{M}_k | \mathbf{Z}_k^{n_k} \right)$ is an $(\mathbf{L}_k, Q_{\Theta_k}, \lambda_k)$ -Gibbs conditional probability measure.

Proof:

From lemma 16, it follows that

$$\begin{aligned} \mathbf{G}_k \left(P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) &= \int \lambda_k \left(D \left(P_{\Theta_k}^{(m)} \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right. \\ &\left. - D \left(P_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel P_{\Theta_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) + D \left(P_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel Q_{\Theta_k} \right) - D \left(P_{\Theta_k}^{(m)} \parallel Q_{\Theta_k} \right) \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0). \end{aligned} \quad (198)$$

The proof proceeds by focusing on the expression $\int D \left(P_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel Q_{\Theta_k} \right) - D \left(P_{\Theta_k}^{(m)} \parallel Q_{\Theta_k} \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0)$, from which it holds that

$$\begin{aligned} &\int D \left(P_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel Q_{\Theta_k} \right) - D \left(P_{\Theta_k}^{(m)} \parallel Q_{\Theta_k} \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0) \\ &= \int D \left(P_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel Q_{\Theta_k} \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0) - D \left(P_{\Theta_k}^{(m)} \parallel Q_{\Theta_k} \right) \end{aligned} \quad (199)$$

$$\begin{aligned} &= \iint \log \left(\frac{dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k)}{dQ_{\Theta_k}} \right) dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) \\ &\quad - \int \log \left(\frac{dP_{\Theta_k}^{(m)}(\theta_k)}{dQ_{\Theta_k}} \right) dP_{\Theta_k}^{(m)}(\theta_k) \end{aligned} \quad (200)$$

$$\begin{aligned} &= \iint \log \left(\frac{dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k)}{dQ_{\Theta_k}} \right) dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) \\ &\quad - \int \log \left(\frac{dP_{\Theta_k}^{(m)}(\theta_k)}{dQ_{\Theta_k}} \right) \left(\int \frac{dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k)}{dP_{\Theta_k}^{(m)}}(\theta_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) \right) dP_{\Theta_k}^{(m)}(\theta_k) \end{aligned} \quad (201)$$

$$= \iint \log \left(\frac{dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k)}{dQ_{\Theta_k}} \right) dP_{\Theta_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\theta_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0)$$

$$- \iint \log \left(\frac{dP_{\Theta_k}^{(m)}}{dQ_{\Theta_k}} (\theta_k) \right) \frac{dP_{\Theta_k|Z_0=z_0}^{(m)} (\theta_k)}{dP_{\Theta_k}^{(m)} (\theta_k)} dP_{\Theta_k}^{(m)} (\theta_k) dP_{Z_0} (z_0) \quad (202)$$

$$= \iint \left(\log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)} (\theta_k)}{dQ_{\Theta_k}} \right) - \log \left(\frac{dP_{\Theta_k}^{(m)} (\theta_k)}{dQ_{\Theta_k}} \right) \right) dP_{\Theta_k|Z_0=z_0}^{(m)} (\theta_k) dP_{Z_0} (z_0) \quad (203)$$

$$= \iint \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)} (\theta_k)}{dQ_{\Theta_k}} \frac{dQ_{\Theta_k}}{dP_{\Theta_k}^{(m)}} (\theta_k) \right) dP_{\Theta_k|Z_0=z_0}^{(m)} (\theta_k) dP_{Z_0} (z_0) \quad (204)$$

$$= \iint \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)} (\theta_k)}{dP_{\Theta_k}^{(m)}} (\theta_k) \right) dP_{\Theta_k|Z_0=z_0}^{(m)} (\theta_k) dP_{Z_0} (z_0) \quad (205)$$

$$= I \left(P_{\Theta_k|Z_0=z_0}^{(m)} ; P_{Z_0} \right), \quad (206)$$

where the equality in (201) follows from [58, Theorem 10]; (202) follows from [67, Theorem 2.6.6]; (203) follows from a change of measure [58, Theorem 2]; (204) follows from [58, Theorem 5]; and (205) follows from [58, Theorem 4]. This completes the proof. \blacksquare

The proof proceeds by focusing on the term $\int D \left(P_{\Theta_k}^{(m)} \parallel P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) dP_{Z_0} (z_0)$ in (197), from which it follows that

$$\begin{aligned} & \int D \left(P_{\Theta_k}^{(m)} \parallel P_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) dP_{Z_0} (z_0) \\ &= \iint \log \left(\frac{dP_{\Theta_k}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\Theta_k}^{(m)} (\theta_k) dP_{Z_0} (z_0) \end{aligned} \quad (207)$$

$$= \iint \log \left(\frac{dP_{\Theta_k}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \frac{dP_{\Theta_k|Z_0=z_0}^{(m)} (\theta_k)}{dP_{\Theta_k|Z_0=z_0}^{(m)}} (\theta_k) \right) dP_{\Theta_k}^{(m)} (\theta_k) dP_{Z_0} (z_0) \quad (208)$$

$$\begin{aligned} &= \iint \log \left(\frac{dP_{\Theta_k}^{(m)}}{dP_{\Theta_k|Z_0=z_0}^{(m)}} (\theta_k) \right) dP_{\Theta_k}^{(m)} (\theta_k) dP_{Z_0} (z_0) \\ &+ \iint \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\Theta_k}^{(m)} (\theta_k) dP_{Z_0} (z_0) \end{aligned} \quad (209)$$

$$\begin{aligned} &= \int D \left(P_{\Theta_k}^{(m)} \parallel P_{\Theta_k|Z_0=z_0}^{(m)} \right) dP_{Z_0} (z_0) \\ &+ \iint \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\Theta_k}^{(m)} (\theta_k) dP_{Z_0} (z_0) \end{aligned} \quad (210)$$

$$= L \left(P_{\Theta_k|Z_0=z_0}^{(m)} ; P_{Z_0} \right) + \iint \log \left(\frac{dP_{\Theta_k|Z_0=z_0}^{(m)}}{dP_{\Theta_k|Z_k=z_k}^{(Q_{\Theta_k}, \lambda_k)}} (\theta_k) \right) dP_{\Theta_k}^{(m)} (\theta_k) dP_{Z_0} (z_0), \quad (211)$$

where the equality in (209) follows from [58, Theorem 8].

From Lemma 26, by using (211) in (197), it holds that

$$\begin{aligned} & \mathbf{G}_k \left(P_{\underline{\boldsymbol{\theta}}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) \\ &= \lambda_k I \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0}^{(m)}; P_{\mathbf{Z}_0} \right) + \lambda_k \left(L \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}; P_{\mathbf{Z}_0} \right) \right. \\ & \quad + \iint \log \left(\frac{dP_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}(\boldsymbol{\theta}_k)}{dP_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)}}(\boldsymbol{\theta}_k)} \right) dP_{\boldsymbol{\theta}_k}^{(m)}(\boldsymbol{\theta}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0) \\ & \quad \left. - \int D \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) dP_{\mathbf{Z}_0}(\mathbf{z}_0) \right) \end{aligned} \quad (212)$$

$$= \lambda_k \left(I \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0}^{(m)}; P_{\mathbf{Z}_0} \right) + L \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0}^{(m)}; P_{\mathbf{Z}_0} \right) + J_{L_k, Q_{\boldsymbol{\theta}_k}, \lambda_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0}^{(m)}; P_{\mathbf{Z}_0} \right) \right), \quad (213)$$

where (213) follows from the definition of the information measure $J_{L_k, Q_{\boldsymbol{\theta}_k}, \lambda_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0}^{(m)}; P_{\mathbf{Z}_0} \right)$ in (92). This completes the proof.

L Proof of Theorem 20

From the definition of the generalization error in (62), it follows that,

$$\begin{aligned} & \mathbf{G}_k \left(P_{\underline{\boldsymbol{\theta}}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) \\ &= \iint \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) - \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) dP_{\mathbf{Z}_k}(\widehat{\mathbf{z}}_k) dP_{\mathbf{Z}_0}(\mathbf{z}_0). \end{aligned} \quad (214)$$

The proof proceeds by focusing on $\mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) - \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right)$, from which it follows that

$$\begin{aligned} & \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) - \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) \\ &= \mathbf{G}_k \left(\widehat{\mathbf{z}}_k, P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}, P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \widehat{\mathbf{z}}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) + \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \widehat{\mathbf{z}}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) \\ & \quad - \mathbf{G}_k \left(\mathbf{z}_k, P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)}, P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) - \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) \end{aligned} \quad (215)$$

$$\begin{aligned} &= \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) - \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \widehat{\mathbf{z}}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) + \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \widehat{\mathbf{z}}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) \\ & \quad - \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \right) + \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) - \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) \end{aligned} \quad (216)$$

$$\begin{aligned} &= \mathbf{R}_{k, \widehat{\mathbf{z}}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \widehat{\mathbf{z}}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) - \mathbf{R}_{k, \mathbf{z}_k} \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) \\ & \quad + \lambda_k \left(D \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \widehat{\mathbf{z}}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \parallel Q_{\boldsymbol{\theta}_k} \right) + D \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \widehat{\mathbf{z}}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) - D \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel Q_{\boldsymbol{\theta}_k} \right) \right) \\ & \quad - \lambda_k \left(D \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \parallel Q_{\boldsymbol{\theta}_k} \right) + D \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel P_{\boldsymbol{\theta}_k | \mathbf{Z}_k = \mathbf{z}_k}^{(Q_{\boldsymbol{\theta}_k}, \lambda_k)} \right) - D \left(P_{\boldsymbol{\theta}_k | \mathbf{Z}_0 = \mathbf{z}_0}^{(m)} \parallel Q_{\boldsymbol{\theta}_k} \right) \right) \end{aligned} \quad (217)$$

$$\begin{aligned}
&= R_{k, \widehat{z}_k} \left(P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) + \lambda_k D \left(P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \parallel Q_{\Theta_k} \right) \\
&\quad - R_{k, z_k} \left(P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) - \lambda_k D \left(P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \parallel Q_{\Theta_k} \right) \\
&\quad + \lambda_k \left(D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) - D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right) \tag{218}
\end{aligned}$$

$$\begin{aligned}
&= R_{k, \widehat{z}_k} (Q_{\Theta_k}) + \lambda_k D \left(Q_{\Theta_k} \parallel P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) - R_{k, z_k} (Q_{\Theta_k}) + \lambda_k D \left(Q_{\Theta_k} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \\
&\quad + \lambda_k \left(D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) - D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right), \tag{219}
\end{aligned}$$

where the equality in (217) follows from Lemma 12; and the equality in (219) follows from Lemma 7. By plugging (219) into (214) it follows that

$$\begin{aligned}
\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{Y}^{(m)} | Z_0}; P_{Z_0} \right) &= \iint R_{k, \widehat{z}_k} (Q_{\Theta_k}) + \lambda_k D \left(Q_{\Theta_k} \parallel P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \\
&\quad - R_{k, z_k} (Q_{\Theta_k}) + \lambda_k D \left(Q_{\Theta_k} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \\
&\quad + \lambda_k \left(D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right. \\
&\quad \left. - D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right) dP_{Z_k}(\widehat{z}_k) dP_{Z_0}(z_0) \tag{220}
\end{aligned}$$

$$\begin{aligned}
&= \iint R_{k, \widehat{z}_k} (Q_{\Theta_k}) + \lambda_k D \left(Q_{\Theta_k} \parallel P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) dP_{Z_k}(\widehat{z}_k) dP_{Z_0}(z_0) \\
&\quad - \iint R_{k, z_k} (Q_{\Theta_k}) + \lambda_k D \left(Q_{\Theta_k} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) dP_{Z_k}(\widehat{z}_k) dP_{Z_0}(z_0) \\
&\quad + \iint \lambda_k \left(D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right. \\
&\quad \left. - D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right) dP_{Z_k}(\widehat{z}_k) dP_{Z_0}(z_0) \tag{221}
\end{aligned}$$

$$\begin{aligned}
&= \lambda_k \iint \left(D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right. \\
&\quad \left. - D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = z_k}^{(Q_{\Theta_k}, \lambda_k)} \right) \right) dP_{Z_k}(\widehat{z}_k) dP_{Z_0}(z_0). \tag{222}
\end{aligned}$$

This completes the proof.

M Proof of Theorem 21

From Theorem 20, it holds that

$$\iint D \left(P_{\Theta_k | Z_0 = z_0}^{(m)} \parallel P_{\Theta_k | Z_k = \widehat{z}_k}^{(Q_{\Theta_k}, \lambda_k)} \right) dP_{Z_k}(\widehat{z}_k) dP_{Z_0}(z_0)$$

$$= \frac{1}{\lambda_k} \mathbf{G}_k \left(P_{\underline{\theta}, \underline{u}^{(m)}, \underline{v}^{(m)} | Z_0}; P_{Z_0} \right) + \int D \left(P_{\underline{\theta}_k | Z_0 = z_0}^{(m)} \parallel P_{\underline{\theta}_k | Z_k = \tilde{z}_k}^{(Q_{\underline{\theta}_k, \lambda_k})} \right) dP_{Z_0}(z_0). \quad (223)$$

Moreover, the term on the left side of the equality in (223) can be expressed as follows

$$\begin{aligned} & \iint D \left(P_{\underline{\theta}_k | Z_0 = z_0}^{(m)} \parallel P_{\underline{\theta}_k | Z_k = \tilde{z}_k}^{(Q_{\underline{\theta}_k, \lambda_k})} \right) dP_{Z_0}(z_0) dP_{Z_k}(\tilde{z}_k) \\ &= \iiint \log \left(\frac{dP_{\underline{\theta}_k | Z_0 = z_0}^{(m)}(\theta_k)}{dP_{\underline{\theta}_k | Z_k = \tilde{z}_k}^{(Q_{\underline{\theta}_k, \lambda_k})}(\theta_k)} \right) dP_{\underline{\theta}_k | Z_0 = z_0}^{(m)}(\theta_k) dP_{Z_0}(z_0) dP_{Z_k}(\tilde{z}_k) \end{aligned} \quad (224)$$

$$= \iiint \log \left(\frac{dP_{\underline{\theta}_k | Z_0 = z_0}^{(m)}(\theta_k)}{dP_{\underline{\theta}_k | Z_k = \tilde{z}_k}^{(Q_{\underline{\theta}_k, \lambda_k})}(\theta_k)} \frac{dP_{\underline{\theta}_k}^{(m)}(\theta_k)}{dP_{\underline{\theta}_k}^{(m)}(\theta_k)} \right) dP_{\underline{\theta}_k | Z_0 = z_0}^{(m)}(\theta_k) dP_{Z_0}(z_0) dP_{Z_k}(\tilde{z}_k) \quad (225)$$

$$\begin{aligned} &= \iiint \log \left(\frac{dP_{\underline{\theta}_k}^{(m)}(\theta_k)}{dP_{\underline{\theta}_k | Z_k = \tilde{z}_k}^{(Q_{\underline{\theta}_k, \lambda_k})}(\theta_k)} \right) dP_{\underline{\theta}_k | Z_0 = z_0}^{(m)}(\theta_k) dP_{Z_0}(z_0) dP_{Z_k}(\tilde{z}_k) \\ &+ \iiint \log \left(\frac{dP_{\underline{\theta}_k | Z_0 = z_0}^{(m)}(\theta_k)}{dP_{\underline{\theta}_k}^{(m)}(\theta_k)} \right) dP_{\underline{\theta}_k | Z_0 = z_0}^{(m)}(\theta_k) dP_{Z_0}(z_0) dP_{Z_k}(\tilde{z}_k) \end{aligned} \quad (226)$$

$$\begin{aligned} &= \iiint \log \left(\frac{dP_{\underline{\theta}_k}^{(m)}(\theta_k)}{dP_{\underline{\theta}_k | Z_k = \tilde{z}_k}^{(Q_{\underline{\theta}_k, \lambda_k})}(\theta_k)} \right) \frac{dP_{\underline{\theta}_k | Z_0 = z_0}^{(m)}(\theta_k)}{dP_{\underline{\theta}_k}^{(m)}(\theta_k)} dP_{\underline{\theta}_k}^{(m)}(\theta_k) dP_{Z_0}(z_0) dP_{Z_k}(\tilde{z}_k) \\ &+ \int D \left(P_{\underline{\theta}_k | Z_0 = z_0}^{(m)} \parallel P_{\underline{\theta}_k}^{(m)} \right) dP_{Z_0}(z_0) \end{aligned} \quad (227)$$

$$\begin{aligned} &= \iint \log \left(\frac{dP_{\underline{\theta}_k}^{(m)}(\theta_k)}{dP_{\underline{\theta}_k | Z_k = \tilde{z}_k}^{(Q_{\underline{\theta}_k, \lambda_k})}(\theta_k)} \right) \left(\int \frac{dP_{\underline{\theta}_k | Z_0 = z_0}^{(m)}(\theta_k)}{dP_{\underline{\theta}_k}^{(m)}(\theta_k)} dP_{Z_0}(z_0) \right) dP_{\underline{\theta}_k}^{(m)}(\theta_k) dP_{Z_k}(\tilde{z}_k) \\ &+ \int D \left(P_{\underline{\theta}_k | Z_0 = z_0}^{(m)} \parallel P_{\underline{\theta}_k}^{(m)} \right) dP_{Z_0}(z_0) \end{aligned} \quad (228)$$

$$= \int D \left(P_{\underline{\theta}_k}^{(m)} \parallel P_{\underline{\theta}_k | Z_k = \tilde{z}_k}^{(Q_{\underline{\theta}_k, \lambda_k})} \right) dP_{Z_k}(\tilde{z}_k) + \int D \left(P_{\underline{\theta}_k | Z_0 = z_0}^{(m)} \parallel P_{\underline{\theta}_k}^{(m)} \right) dP_{Z_0}(z_0), \quad (229)$$

where (225) follows from [58, Theorem 3]; (226) follows from [58, Theorem 4]; (227) follows from [58, Theorem 2]; and (228) follows from [67, Theorem 2.6.6] and [58, Theorem 10].

Using (229) in (223) yields

$$\begin{aligned} & \int D \left(P_{\underline{\theta}_k}^{(m)} \parallel P_{\underline{\theta}_k | Z_k = \tilde{z}_k}^{(Q_{\underline{\theta}_k, \lambda_k})} \right) dP_{Z_k}(\tilde{z}_k) + \int D \left(P_{\underline{\theta}_k | Z_0 = z_0}^{(m)} \parallel P_{\underline{\theta}_k}^{(m)} \right) dP_{Z_0}(z_0) \\ &= \frac{1}{\lambda_k} \mathbf{G}_k \left(P_{\underline{\theta}, \underline{u}^{(m)}, \underline{v}^{(m)} | Z_0}; P_{Z_0} \right) + \int D \left(P_{\underline{\theta}_k | Z_0 = z_0}^{(m)} \parallel P_{\underline{\theta}_k | Z_k = \tilde{z}_k}^{(Q_{\underline{\theta}_k, \lambda_k})} \right) dP_{Z_0}(z_0), \end{aligned} \quad (230)$$

which completes the proof.

N Proof of Theorem 22

The proof follows along the same lines as the proof of Theorem 18. It is included entirely hereunder for the sake of completeness.

The proof starts by considering $J_{L_k, Q_{Z_k}, \alpha_k} \left(P_{Z_k | \Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right)$ in (99), from which it follows that

$$\begin{aligned}
& J_{L_k, Q_{Z_k}, \alpha_k} \left(P_{Z_k | \Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) \\
&= \int_{\mathcal{A}_\gamma} \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) dP_{Z_k} P_{\Theta_k}^{(m)}(z_k, \theta_k) \\
&+ \int_{\mathcal{A}_\gamma^c} \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) dP_{Z_k} P_{\Theta_k}^{(m)}(z_k, \theta_k) \\
&- \int_{\mathcal{A}_\gamma} \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) dP_{Z_k | \Theta_k}^{(m)} P_{\Theta_k}^{(m)}(z_k, \theta_k) \\
&- \int_{\mathcal{A}_\gamma^c} \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) dP_{Z_k | \Theta_k}^{(m)} P_{\Theta_k}^{(m)}(z_k, \theta_k) \tag{231}
\end{aligned}$$

$$\begin{aligned}
& \leq \gamma P_{Z_k} P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma^c) - \gamma P_{Z_k | \Theta_k}^{(m)} P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma) \\
&+ \max_{(\theta_k, z_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) P_{Z_k} P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma) \\
&- \min_{(\theta_k, z_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) P_{Z_k | \Theta_k}^{(m)} P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma^c) \tag{232}
\end{aligned}$$

$$\begin{aligned}
& = -\gamma P_{Z_k} P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma) + \gamma P_{Z_k | \Theta_k}^{(m)} P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma^c) \\
&+ \max_{(\theta_k, z_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) P_{Z_k} P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma) \\
&- \min_{(\theta_k, z_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) P_{Z_k | \Theta_k}^{(m)} P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma^c) \tag{233} \\
& = P_{Z_k} P_{\Theta_k}^{(m)}(\mathcal{A}_\gamma) \left(\max_{(\theta_k, z_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{\tilde{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) - \gamma \right)
\end{aligned}$$

$$+P_{\mathbf{Z}_k|\Theta_k}^{(m)} P_{\Theta_k}^{(m)} \left(\mathcal{A}_\gamma^c \right) \left(\gamma - \min_{(\theta_k, \mathbf{z}_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}(\mathbf{z}_k)} \right) \right), \quad (234)$$

where the inequality in (232) follows from the definition of the set \mathcal{A}_γ and the extrema of the log-likelihood ratio.

A lower bounds on the information measure $J_{L_k, Q_{\mathbf{Z}_k}, \alpha_k} \left(P_{\mathbf{Z}_k|\Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right)$ in (99) can also be obtained as follows,

$$\begin{aligned} & J_{L_k, Q_{\Theta_k}, \lambda_k} \left(P_{\Theta_k|Z_0}^{(m)}; P_{Z_0} \right) \\ &= \int_{\mathcal{A}_\gamma} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}(\mathbf{z}_k)} \right) dP_{\Theta_k}^{(m)} P_{Z_0}(\theta_k, \mathbf{z}_0) \\ &+ \int_{\mathcal{A}_\gamma^c} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}(\mathbf{z}_k)} \right) dP_{\Theta_k}^{(m)} P_{Z_0}(\theta_k, \mathbf{z}_0) \\ &- \int_{\mathcal{A}_\gamma} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}(\mathbf{z}_k)} \right) dP_{\Theta_k|Z_0}^{(m)} P_{Z_0}(\theta_k, \mathbf{z}_0) \\ &- \int_{\mathcal{A}_\gamma^c} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}(\mathbf{z}_k)} \right) dP_{\Theta_k|Z_0}^{(m)} P_{Z_0}(\theta_k, \mathbf{z}_0) \end{aligned} \quad (235)$$

$$\begin{aligned} & \geq \gamma P_{\Theta_k}^{(m)} P_{Z_0}(\mathcal{A}_\gamma) - \gamma P_{\mathbf{Z}_k|\Theta_k}^{(m)} P_{\Theta_k}^{(m)} \left(\mathcal{A}_\gamma^c \right) \\ &+ \min_{(\theta_k, \mathbf{z}_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}(\mathbf{z}_k)} \right) P_{\Theta_k}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma^c \right) \\ &- \max_{(\theta_k, \mathbf{z}_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}(\mathbf{z}_k)} \right) P_{\mathbf{Z}_k|\Theta_k}^{(m)} P_{\Theta_k}^{(m)} \left(\mathcal{A}_\gamma \right) \end{aligned} \quad (236)$$

$$\begin{aligned} &= -\gamma P_{\Theta_k}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma^c \right) + \gamma P_{\mathbf{Z}_k|\Theta_k}^{(m)} P_{\Theta_k}^{(m)} \left(\mathcal{A}_\gamma \right) \\ &+ \min_{(\theta_k, \mathbf{z}_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}(\mathbf{z}_k)} \right) P_{\Theta_k}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma^c \right) \\ &- \max_{(\theta_k, \mathbf{z}_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}(\mathbf{z}_k)} \right) P_{\mathbf{Z}_k|\Theta_k}^{(m)} P_{\Theta_k}^{(m)} \left(\mathcal{A}_\gamma \right) \end{aligned} \quad (237)$$

$$\begin{aligned}
&= P_{\Theta_k}^{(m)} P_{Z_0} \left(\mathcal{A}_\gamma^c \right) \left(\min_{(\theta_k, z_0) \in \mathcal{A}_\gamma^c} \log \left(\frac{dP_{Z_k|\Theta_k=\theta_k}^{(m)}}{dP_{\tilde{Z}_k|\Theta_k=\theta_k}^{(Q_{Z_k}, \alpha_k)}} (z_k) \right) - \gamma \right) \\
&\quad + P_{Z_k|\Theta_k}^{(m)} P_{\Theta_k}^{(m)} \left(\mathcal{A}_\gamma \right) \left(\gamma - \max_{(\theta_k, z_0) \in \mathcal{A}_\gamma} \log \left(\frac{dP_{Z_k|\Theta_k=\theta_k}^{(m)}}{dP_{\tilde{Z}_k|\Theta_k=\theta_k}^{(Q_{Z_k}, \alpha_k)}} (z_k) \right) \right), \quad (238)
\end{aligned}$$

where the equality in (236) follows from the definition of the set \mathcal{A}_γ and the extrema of the log-likelihood ratio.

O Proof of Theorem 23

The proof follows along the same lines as the proof of Theorem 19. It is included entirely hereunder for the sake of completeness.

The following lemma is introduced as part of this proof.

Lemma 27. *Consider the generalization error $\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}; P_{Z_0} \right)$ in (62) and assume that for all $\theta_k \in \mathcal{M}_k$, the probability measures $P_{Z_k|\Theta_k=\theta_k}^{(m)}$ in (56); P_{Z_k} in (32); and the σ -finite measure Q_{Z_k} in (75) satisfy $P_{Z_k|\Theta_k=\theta_k}^{(m)} \ll Q_{Z_k} \ll P_{Z_k} \ll Q_{Z_k}$. Then,*

$$\begin{aligned}
\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}; P_{Z_0} \right) &= \alpha_k I \left(P_{Z_k|\Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) \\
&+ \alpha_k \int \left(D \left(P_{Z_k} \parallel P_{\tilde{Z}_k|\Theta_k=\theta_k}^{(Q_{Z_k}, \alpha_k)} \right) - D \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \parallel P_{\tilde{Z}_k|\Theta_k=\theta_k}^{(Q_{Z_k}, \alpha_k)} \right) \right) dP_{\Theta_k}^{(m)}(\theta_k), \quad (239)
\end{aligned}$$

where the conditional probability measure $P_{\tilde{Z}_k|\Theta_k}^{(Q_{Z_k}, \alpha_k)} \in \Delta \left(\mathcal{Z}_k^{n_k} | \mathcal{M}_k \right)$ is an $(\mathbb{L}_k, Q_{Z_k}, \alpha_k)$ -Gibbs conditional probability measure; and the probability measures $P_{\Theta_k}^{(m)}$ is defined in (55).

Proof: From lemma 17 it follows that

$$\begin{aligned}
\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}; P_{Z_0} \right) &= \int \alpha_k \left(D \left(P_{Z_k} \parallel P_{\tilde{Z}_k|\Theta_k=\theta_k}^{(Q_{Z_k}, \alpha_k)} \right) \right. \\
&\left. - D \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \parallel P_{\tilde{Z}_k|\Theta_k=\theta_k}^{(Q_{Z_k}, \alpha_k)} \right) + D \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \parallel Q_{Z_k} \right) - D \left(P_{Z_k} \parallel Q_{Z_k} \right) \right) dP_{\Theta_k}^{(m)}(\theta_k). \quad (240)
\end{aligned}$$

The proof proceeds by focusing on the expression $\int D \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \parallel Q_{Z_k} \right) - D \left(P_{Z_k} \parallel Q_{Z_k} \right) dP_{\Theta_k}^{(m)}(\theta_k)$, from which it holds that

$$\int D \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \parallel Q_{Z_k} \right) - D \left(P_{Z_k} \parallel Q_{Z_k} \right) dP_{\Theta_k}^{(m)}(\theta_k)$$

$$= \int D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel Q_{\mathbf{Z}_k} \right) dP_{\Theta_k}^{(m)}(\theta_k) - D(P_{\mathbf{Z}_k} \parallel Q_{\mathbf{Z}_k}) \quad (241)$$

$$= \iint \log \left(\frac{dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k)}{dQ_{\mathbf{Z}_k}}(\mathbf{z}_k) \right) dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \quad (242)$$

$$\begin{aligned} & - \int \log \left(\frac{dP_{\mathbf{Z}_k}(\mathbf{z}_k)}{dQ_{\mathbf{Z}_k}}(\mathbf{z}_k) \right) dP_{\mathbf{Z}_k}(\mathbf{z}_k) \\ & = \iint \log \left(\frac{dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k)}{dQ_{\mathbf{Z}_k}}(\mathbf{z}_k) \right) dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \\ & - \int \log \left(\frac{dP_{\mathbf{Z}_k}(\mathbf{z}_k)}{dQ_{\mathbf{Z}_k}}(\mathbf{z}_k) \right) \left(\int \frac{dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\mathbf{Z}_k}}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \right) dP_{\mathbf{Z}_k}(\mathbf{z}_k) \quad (243) \end{aligned}$$

$$\begin{aligned} & = \iint \log \left(\frac{dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k)}{dQ_{\mathbf{Z}_k}}(\mathbf{z}_k) \right) dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \\ & - \iint \log \left(\frac{dP_{\mathbf{Z}_k}(\mathbf{z}_k)}{dQ_{\mathbf{Z}_k}}(\mathbf{z}_k) \right) \frac{dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\mathbf{Z}_k}}(\mathbf{z}_k) dP_{\mathbf{Z}_k}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \quad (244) \end{aligned}$$

$$= \iint \left(\log \left(\frac{dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k)}{dQ_{\mathbf{Z}_k}}(\mathbf{z}_k) \right) - \log \left(\frac{dP_{\mathbf{Z}_k}(\mathbf{z}_k)}{dQ_{\mathbf{Z}_k}}(\mathbf{z}_k) \right) \right) dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \quad (245)$$

$$= \iint \log \left(\frac{dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k)}{dQ_{\mathbf{Z}_k}}(\mathbf{z}_k) \frac{dQ_{\mathbf{Z}_k}}{dP_{\mathbf{Z}_k}}(\mathbf{z}_k) \right) dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \quad (246)$$

$$= \iint \log \left(\frac{dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\mathbf{Z}_k}}(\mathbf{z}_k) \right) dP_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \quad (247)$$

$$= I \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right), \quad (248)$$

where the equality (243) follows from [58, Theorem 10]; (244) follows from [67, Theorem 2.6.6]; (245) follows from a change of measure [58, Theorem 2]; (246) follows from [58, Theorem 5]; and (247) follows from [58, Theorem 4]. This completes the proof. \blacksquare

The proof proceeds by focusing on $\int D \left(P_{\mathbf{Z}_k} \parallel P_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)}(\theta_k)$ in (239), from which it follows that,

$$\begin{aligned} & \int D \left(P_{\mathbf{Z}_k} \parallel P_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)}(\theta_k) \\ & = \iint \log \left(\frac{dP_{\mathbf{Z}_k}(\mathbf{z}_k)}{dP_{\tilde{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}(\mathbf{z}_k) \right) dP_{\mathbf{Z}_k}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \quad (249) \end{aligned}$$

$$= \iint \log \left(\frac{dP_{\mathbf{Z}_k}}{dP_{\widehat{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}(\mathbf{z}_k)} \frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)} \right) dP_{\mathbf{Z}_k}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \quad (250)$$

$$\begin{aligned} &= \iint \log \left(\frac{dP_{\mathbf{Z}_k}}{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}}(\mathbf{z}_k) \right) (\mathbf{z}_k) dP_{\mathbf{Z}_k}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \\ &+ \iint \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\widehat{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}(\mathbf{z}_k)} \right) dP_{\mathbf{Z}_k}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \end{aligned} \quad (251)$$

$$= \int D \left(P_{\mathbf{Z}_k} \parallel P_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)} \right) dP_{\Theta_k}^{(m)}(\theta_k) \quad (252)$$

$$\begin{aligned} &+ \iint \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\widehat{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}(\mathbf{z}_k) \right) dP_{\mathbf{Z}_k}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \\ &= L \left(P_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) + \iint \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\widehat{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}(\mathbf{z}_k) \right) dP_{\mathbf{Z}_k}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k), \end{aligned} \quad (253)$$

where the equality in (251) follows from [58, Theorem 8].

From Lemma 27, by using (253) in (239), it holds that

$$\begin{aligned} &\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)}|\mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) \\ &= \alpha_k I \left(P_{\mathbf{Z}_k|\Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) + \alpha_k \left(L \left(P_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) \right. \\ &\quad \left. + \iint \log \left(\frac{dP_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)}(\mathbf{z}_k)}{dP_{\widehat{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}}(\mathbf{z}_k) \right) dP_{\mathbf{Z}_k}(\mathbf{z}_k) dP_{\Theta_k}^{(m)}(\theta_k) \right) \end{aligned} \quad (254)$$

$$- \int D \left(P_{\mathbf{Z}_k|\Theta_k=\theta_k}^{(m)} \parallel P_{\widehat{\mathbf{Z}}_k|\Theta_k=\theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)}(\theta_k) \quad (255)$$

$$= \alpha_k \left(I \left(P_{\mathbf{Z}_k|\Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) + L \left(P_{\mathbf{Z}_k|\Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) + J_{L_k, Q_{\mathbf{Z}_k}, \alpha_k} \left(P_{\mathbf{Z}_k|\Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right) \right), \quad (256)$$

where (256) follows from the definition of the information measure $J_{L_k, Q_{\mathbf{Z}_k}, \alpha_k} \left(P_{\mathbf{Z}_k|\Theta_k}^{(m)}; P_{\Theta_k}^{(m)} \right)$ in (99). This completes the proof.

P Proof of Theorem 24

The proof follows along the same lines as the proof of Theorem 20. It is included entirely hereunder for the sake of completeness.

From lemma 5, it holds that

$$\mathbf{G}_k \left(P_{\underline{\Theta}, \underline{\mathbf{U}}^{(m)}, \underline{\mathbf{V}}^{(m)}|\mathbf{Z}_0}; P_{\mathbf{Z}_0} \right)$$

$$= \iint \left(R_{k, \hat{\theta}_k} \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \right) - R_{k, \theta_k} \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \right) \right) dP_{\Theta_k}(\hat{\theta}_k) dP_{\Theta_k}(\theta_k). \quad (257)$$

The proof proceeds by focusing on $R_{k, \hat{\theta}_k} \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \right) - R_{k, \theta_k} \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \right)$, from which it follows that

$$\begin{aligned} & R_{k, \hat{\theta}_k} \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \right) - R_{k, \theta_k} \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \right) \\ &= R_{k, \hat{\theta}_k} \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) - G_k \left(\hat{\theta}_k, P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}, P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \right) \\ & \quad + G_k \left(\theta_k, P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)}, P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \right) - R_{k, \theta_k} \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \\ &= R_{k, \hat{\theta}_k} \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) - R_{k, \theta_k} \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \end{aligned} \quad (258)$$

$$\begin{aligned} & -\alpha_k \left(-D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) + D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel Q_{\mathbf{Z}_k} \right) \right. \\ & \quad \left. - D \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) \right) + \alpha_k \left(-D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \right. \\ & \quad \left. + D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel Q_{\mathbf{Z}_k} \right) - D \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) \right) \end{aligned} \quad (259)$$

$$\begin{aligned} &= R_{k, \hat{\theta}_k} \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) + \alpha_k D \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) \\ & \quad - \left(R_{k, \theta_k} \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) + \alpha_k D \left(P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \parallel Q_{\mathbf{Z}_k} \right) \right) \\ & \quad + \alpha_k \left(D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) - D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \right) \end{aligned} \quad (260)$$

$$\begin{aligned} &= R_{k, \hat{\theta}_k} \left(Q_{\mathbf{Z}_k} \right) - \alpha_k D \left(Q_{\mathbf{Z}_k} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) - \left(R_{k, \theta_k} \left(Q_{\mathbf{Z}_k} \right) - \alpha_k D \left(Q_{\mathbf{Z}_k} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \right) \\ & \quad + \alpha_k \left(D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) - D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \right), \end{aligned} \quad (261)$$

where the equality in (259) follows from Lemma 13; and the equality in (261) follows from Lemma 9. By plugging (261) into (257), it follows that

$$\begin{aligned} & G_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | \mathbf{Z}_0}; P_{\mathbf{Z}_0} \right) = \iint R_{k, \hat{\theta}_k} \left(Q_{\mathbf{Z}_k} \right) - \alpha_k D \left(Q_{\mathbf{Z}_k} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \\ & \quad - \left(R_{k, \theta_k} \left(Q_{\mathbf{Z}_k} \right) - \alpha_k D \left(Q_{\mathbf{Z}_k} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \right) + \alpha_k \left(D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \right. \\ & \quad \left. - D \left(P_{\mathbf{Z}_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \theta_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) \right) dP_{\Theta_k}(\hat{\theta}_k) dP_{\Theta_k}(\theta_k) \\ &= \iint R_{k, \hat{\theta}_k} \left(Q_{\mathbf{Z}_k} \right) - \alpha_k D \left(Q_{\mathbf{Z}_k} \parallel P_{\hat{\mathbf{Z}}_k | \Theta_k = \hat{\theta}_k}^{(Q_{\mathbf{Z}_k}, \alpha_k)} \right) dP_{\Theta_k}(\hat{\theta}_k) dP_{\Theta_k}(\theta_k) \end{aligned} \quad (262)$$

$$\begin{aligned}
& - \iint \left(R_{k, \theta_k} (Q_{Z_k}) - \alpha_k D \left(Q_{Z_k} \parallel P_{\widehat{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)} \right) \right) dP_{\Theta_k} (\widehat{\theta}_k) dP_{\Theta_k} (\theta_k) \\
& + \iint \alpha_k \left(D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\widehat{Z}_k | \Theta_k = \widehat{\theta}_k}^{(Q_{Z_k}, \alpha_k)} \right) \right. \\
& \left. - D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\widehat{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)} \right) \right) dP_{\Theta_k} (\widehat{\theta}_k) dP_{\Theta_k} (\theta_k) \tag{263}
\end{aligned}$$

$$\begin{aligned}
& = \iint \alpha_k \left(D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\widehat{Z}_k | \Theta_k = \widehat{\theta}_k}^{(Q_{Z_k}, \alpha_k)} \right) \right. \\
& \left. - D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\widehat{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)} \right) \right) dP_{\Theta_k} (\widehat{\theta}_k) dP_{\Theta_k} (\theta_k). \tag{264}
\end{aligned}$$

This completes the proof.

Q Proof of Theorem 25

The proof follows along the same lines as the proof of Theorem 21. It is included entirely hereunder for the sake of completeness.

From Theorem 24, it holds that

$$\begin{aligned}
& \iint D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\widehat{Z}_k | \Theta_k = \widehat{\theta}_k}^{(Q_{Z_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)} (\widehat{\theta}_k) dP_{\Theta_k}^{(m)} (\theta_k) \\
& = \frac{1}{\alpha_k} \mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)} | Z_0}; P_{Z_0} \right) + \int D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\widehat{Z}_k | \Theta_k = \theta_k}^{(Q_{Z_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)} (\theta_k). \tag{265}
\end{aligned}$$

Moreover, the term on the left of the equality in (265) can be expressed as follows

$$\begin{aligned}
& \iint D \left(P_{Z_k | \Theta_k = \theta_k}^{(m)} \parallel P_{\widehat{Z}_k | \Theta_k = \widehat{\theta}_k}^{(Q_{Z_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)} (\widehat{\theta}_k) dP_{\Theta_k}^{(m)} (\theta_k) \\
& = \iiint \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{\widehat{Z}_k | \Theta_k = \widehat{\theta}_k}^{(Q_{Z_k}, \alpha_k)}} \right) dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k) dP_{\Theta_k}^{(m)} (\widehat{\theta}_k) dP_{\Theta_k}^{(m)} (\theta_k) \tag{266}
\end{aligned}$$

$$\begin{aligned}
& = \iiint \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k) dP_{Z_k}(z_k)}{dP_{\widehat{Z}_k | \Theta_k = \widehat{\theta}_k}^{(Q_{Z_k}, \alpha_k)}} \right) dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k) dP_{\Theta_k}^{(m)} (\widehat{\theta}_k) dP_{\Theta_k}^{(m)} (\theta_k) \tag{267}
\end{aligned}$$

$$\begin{aligned}
& = \iiint \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k) dP_{Z_k}(z_k)}{dP_{\widehat{Z}_k | \Theta_k = \widehat{\theta}_k}^{(Q_{Z_k}, \alpha_k)}} \right) dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k) dP_{\Theta_k}^{(m)} (\widehat{\theta}_k) dP_{\Theta_k}^{(m)} (\theta_k) \tag{268}
\end{aligned}$$

$$\begin{aligned}
& = \iint \log \left(\frac{dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k)}{dP_{Z_k}} \right) dP_{Z_k | \Theta_k = \theta_k}^{(m)}(z_k) dP_{\Theta_k}^{(m)} (\theta_k)
\end{aligned}$$

$$+ \iiint \log \left(\frac{dP_{Z_k}}{dP_{\tilde{Z}_k|\Theta_k=\hat{\theta}_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) dP_{Z_k|\Theta_k=\theta_k}^{(m)}(z_k) dP_{\Theta_k}^{(m)}(\hat{\theta}_k) dP_{\Theta_k}^{(m)}(\theta_k) \quad (269)$$

$$= \int D \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \parallel P_{Z_k} \right) dP_{\Theta_k}^{(m)}(\theta_k)$$

$$+ \iiint \log \left(\frac{dP_{Z_k}}{dP_{\tilde{Z}_k|\Theta_k=\hat{\theta}_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) \frac{dP_{Z_k|\Theta_k=\theta_k}^{(m)}(z_k)}{dP_{Z_k}} dP_{Z_k}(z_k) dP_{\Theta_k}^{(m)}(\hat{\theta}_k) dP_{\Theta_k}^{(m)}(\theta_k) \quad (270)$$

$$= \int D \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \parallel P_{Z_k} \right) dP_{\Theta_k}^{(m)}(\theta_k)$$

$$+ \iint \log \left(\frac{dP_{Z_k}}{dP_{\tilde{Z}_k|\Theta_k=\hat{\theta}_k}^{(Q_{Z_k}, \alpha_k)}}(z_k) \right) \left(\int \frac{dP_{Z_k|\Theta_k=\theta_k}^{(m)}(z_k)}{dP_{Z_k}} dP_{\Theta_k}^{(m)}(\theta_k) \right) dP_{Z_k}(z_k) dP_{\Theta_k}^{(m)}(\hat{\theta}_k) \quad (271)$$

$$= \int D \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \parallel P_{Z_k} \right) dP_{\Theta_k}^{(m)}(\theta_k) + \int D \left(P_{Z_k} \parallel P_{\tilde{Z}_k|\Theta_k=\hat{\theta}_k}^{(Q_{Z_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)}(\hat{\theta}_k), \quad (272)$$

where (267) follows from [58, Theorem 3]; in (268) follows from [58, Theorem 4]; (270) follows from [58, Theorem 2]; and (272) follows from [67, Theorem 2.6.6] and [58, Theorem 10].

Using (272) in (265) yields

$$\begin{aligned} & \int D \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \parallel P_{Z_k} \right) dP_{\Theta_k}^{(m)}(\theta_k) + \int D \left(P_{Z_k} \parallel P_{\tilde{Z}_k|\Theta_k=\hat{\theta}_k}^{(Q_{Z_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)}(\hat{\theta}_k) \\ &= \frac{1}{\alpha_k} \mathbf{G}_k \left(P_{\underline{\Theta}, \underline{U}^{(m)}, \underline{V}^{(m)}|Z_0}; P_{Z_0} \right) + \int D \left(P_{Z_k|\Theta_k=\theta_k}^{(m)} \parallel P_{\tilde{Z}_k|\Theta_k=\theta_k}^{(Q_{Z_k}, \alpha_k)} \right) dP_{\Theta_k}^{(m)}(\theta_k), \quad (273) \end{aligned}$$

which completes the proof.



**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