

Re: Resubmission of manuscript *Architecture-Aware Generalization Bounds for Temporal Networks: Theory and Fair Comparison Methodology* (Manuscript ID: TMLR 5579)

We are pleased to resubmit our manuscript for consideration in *TMLR*. This version is a substantial revision of the previous submission and incorporates the changes promised in our rebuttal.

In brief, the main updates are:

- **Corrected theoretical results.** Lemma 2 now explicitly follows the architecture-aware Rademacher bound of Golowich et al., with the full product of layer spectral norms $R = \prod_{\ell=1}^D M^{(\ell)}$ and \sqrt{D} depth dependence. The main theorem has been restated to clearly separate the imported i.i.d. class bound from our new mixing-based framework for β -mixing sequences.
- **Online-to-batch and convexity.** Proposition 1 has been revised to make the convex Lipschitz loss assumption explicit and to clearly invoke Jensen’s inequality in the online-to-batch step. The resulting generalization bound depends on the number of blocks B , yielding a concentration rate of order $\sqrt{\log N/\bar{N}}$ under exponentially decaying β -mixing.
- **Scope and positioning.** Throughout the paper we now emphasize that we do *not* extend or modify the core theorem of Golowich et al. Instead, we extend the *applicability* of such architecture-aware bounds from i.i.d. samples to β -mixing time series by embedding their class bound in a blocking and delayed-feedback framework.
- **Experiments and presentation.** We clarified the “fair comparison” protocol, softened statements that might over-interpret the bounds (e.g., depth vs. data requirements), and streamlined the exposition and figures to better align the empirical results with the revised theory.

A detailed, point-by-point response to all reviewer and Associate Editor comments is provided in the accompanying rebuttal. We believe these changes fully address the concerns raised in the previous round and significantly strengthen the manuscript.

Thank you very much for your time and consideration.

Architecture-Aware Generalization Bounds for Temporal Networks: Theory and Fair Comparison Methodology

Anonymous authors

Paper under double-blind review

Abstract

Deep temporal architectures such as Temporal Convolutional Networks (TCNs) achieve strong predictive performance on sequential data, yet theoretical understanding of their generalization remains limited. We address this gap through three contributions: introducing a principled evaluation methodology for temporal models, revealing surprising empirical phenomena about temporal dependence, and establishing the first architecture-aware theoretical framework for dependent sequences.

Fair-Comparison Methodology. We introduce evaluation protocols that fix effective sample size N_{eff} to isolate temporal structure effects from information content. This addresses a fundamental challenge: temporal dependence affects both information content and learning dynamics, and standard evaluations conflate these effects. Our methodology enables principled comparison of models across dependency regimes.

Empirical Findings. Applying this methodology reveals that under controlled $N_{\text{eff}} = 2,000$, strongly dependent sequences ($\rho = 0.8$) exhibit approximately 76% smaller generalization gaps than weakly dependent ones ($\rho = 0.2$), challenging the conventional view that dependence universally impedes learning. However, observed convergence rates ($N_{\text{eff}}^{-1.21}$ to $N_{\text{eff}}^{-0.89}$) significantly exceed theoretical worst-case predictions ($N^{-0.5}$), revealing that temporal architectures exploit problem structure in ways current theory does not capture.

Theoretical Framework. To provide the foundations for these empirical investigations, we develop the first architecture-aware generalization bounds for deep temporal models on exponentially β -mixing sequences. By embedding Golowich et al.’s i.i.d. class bound within a novel blocking scheme that partitions N samples into approximately $B \approx N/\log N$ quasi-independent blocks, we establish polynomial sample complexity under convex Lipschitz losses. The framework achieves \sqrt{D} depth scaling alongside the product of layer-wise norms $R = \prod_{\ell=1}^D M^{(\ell)}$, avoiding exponential dependence. While these bounds are conservative, as our empirical results demonstrate, they prove learnability and identify architectural scaling laws, providing worst-case baselines that highlight where future theory must improve to explain observed performance.

1 Introduction

Modern deep architectures, notably Temporal Convolutional Networks (TCNs) Lea et al. (2017); Bai et al. (2018) and Transformer variants Vaswani et al. (2017), underpin state-of-the-art forecasting and representation learning across domains ranging from intensive care monitoring to backbone network management Lim et al. (2021); Oreshkin et al. (2019). Despite this empirical success, two fundamental gaps limit our understanding of temporal deep learning. First, we lack proper evaluation methodology for dependent data that separates temporal structure effects from information content—standard approaches vary sequence length without controlling for effective sample size, confounding temporal structure with statistical information density. Second, we lack theoretical guarantees that explicitly account for architectural choices in temporal models on dependent sequences. These gaps leave researchers without principled answers to practical questions: how should models be compared across different dependency structures, when do dependencies help versus hinder learning, and how deep should temporal networks be?

Understanding temporal model performance requires both evaluation methodology and theoretical foundations. On the evaluation front, standard approaches compare models by varying raw sequence length N without controlling for effective sample size N_{eff} , which accounts for information loss due to temporal correlation. This conflates two distinct effects: changes in temporal structure (correlation strength) and changes in information content (equivalent independent samples). The result is systematic bias where conclusions about "how dependence affects learning" actually reflect confounded comparisons of unequal information budgets.

On the theoretical front, classical Probably Approximately Correct (PAC) theory Valiant (1984) presumes independent observations, making its bounds vacuous for time-series data. Extensions to dependent settings through mixing coefficients Yu (1994); Kuznetsov & Mohri (2017) and sequential Rademacher analyses Rakhlin et al. (2010); Chen et al. (2021) either explode exponentially with depth or depend on norms that grow during training, providing little practical guidance for modern deep temporal architectures.

Our theoretical framework builds on Golowich et al. (2018), who established the first architecture-aware generalization bounds for deep networks on i.i.d. data, showing that Rademacher complexity scales as $\mathcal{O}(\prod_{\ell=1}^D M^{(\ell)} \sqrt{D}/\sqrt{N})$ for networks with depth D and layer-wise spectral norm bounds $M^{(\ell)}$. This achieves \sqrt{D} scaling alongside a product of layer norms—a dramatic improvement over naive VC-dimension analysis yielding exponential depth dependence. Rather than extending or modifying their core complexity analysis, we import their i.i.d. class bound and embed it within a principled framework for dependent sequences via blocking, coupling through Bradley’s inequality, and block-level concentration under β -mixing. Thus, our contribution lies in developing the first complete pipeline from architecture-aware class complexity to generalization guarantees for temporal data, not in improving the underlying architectural analysis itself.

We address both gaps through three complementary contributions spanning evaluation methodology, empirical analysis, and theoretical foundations.

(1) Fair-Comparison Methodology. We introduce evaluation protocols that fix effective sample size N_{eff} when comparing models across dependency structures, isolating temporal structure effects from information content. This addresses a fundamental evaluation challenge: when comparing a sequence with $\rho = 0.2$ (weakly dependent) to one with $\rho = 0.8$ (strongly dependent), should they have the same raw length N or the same information content N_{eff} ? Standard approaches use the former, confounding structural effects with information differences. Our methodology enables principled comparison by equalizing N_{eff} , revealing phenomena invisible to conventional evaluation.

(2) Empirical Findings on Dependence Benefits. Applying this methodology reveals surprising results: at fixed $N_{\text{eff}} = 2,000$, strongly dependent sequences ($\rho = 0.8$) achieve approximately 76% smaller generalization gaps than weakly dependent ones ($\rho = 0.2$), with statistical significance ($p < 0.001$, $n = 12$ trials). This challenges the conventional view that temporal dependence universally impedes learning, suggesting instead that architectural inductive biases can exploit temporal regularities. However, observed convergence rates ($N_{\text{eff}}^{-1.21}$ to $N_{\text{eff}}^{-0.89}$) substantially exceed theoretical worst-case predictions ($N^{-0.5}$), indicating that temporal architectures exploit problem structure in ways current theory does not capture.

(3) Architecture-Aware Theoretical Framework. To provide rigorous foundations for these investigations, we establish the first architecture-aware generalization bounds for deep temporal models on exponentially β -mixing sequences. By embedding Golowich et al.’s i.i.d. class bound within a novel blocking scheme, we prove polynomial sample complexity under convex Lipschitz losses. The framework achieves \sqrt{D} depth scaling alongside the product of layer-wise norms $R = \prod_{\ell=1}^D M^{(\ell)}$, avoiding exponential dependence on depth. Our delayed-feedback blocking mechanism partitions N samples into approximately $B \approx N/\log N$ quasi-independent blocks, incurring only a mild $\sqrt{\log N}$ penalty compared to i.i.d. rates. While these bounds are conservative, as our empirical results demonstrate, they establish learnability and identify architectural scaling laws, providing worst-case baselines that highlight where future theory must improve.

Terminology Note: Throughout this paper, we use “effective sample size” (denoted N_{eff}) to refer to the equivalent number of independent observations that would provide the same statistical information as a

dependent sequence of length N . We distinguish between “standard evaluation” (varying raw sequence length N) and “fair comparison evaluation” (controlling for effective sample size N_{eff}).

Together, these contributions shift our understanding of temporal deep learning in three ways. First, they provide researchers with tools for evaluating temporal models that avoid systematic biases from confounded comparisons. Second, they reveal that temporal dependence can be beneficial rather than detrimental under appropriate architectural design—a finding that challenges conventional wisdom and opens new research directions. Third, they establish a theoretical foundations proving that deep temporal architectures remain viable under dependence with polynomial sample complexity, while the theory-practice gap precisely identifies where future work must develop problem-dependent complexity measures beyond worst-case β -mixing analysis.

The remainder of this paper is structured as follows: Section 2 reviews related work in generalization theory for dependent data and deep learning. Section 3 introduces essential preliminaries, including β -mixing processes, Rademacher complexity, and PAC learning theory. Section 4 presents our architectural generalization bounds for temporal models under β -mixing. Section 5 empirically validates these bounds using synthetic and physiological time series, while Section 6 discusses the implications of our findings for temporal model design. Section 7 concludes with a summary of contributions.

2 Related Work

PAC Learning under Dependence. Understanding how temporal dependencies affect learning guarantees has been a research challenge in machine learning (ML) theory. This area of inquiry began with Yu’s work Yu (1994), which established concentration inequalities for mixing processes, mathematical tools for bounding the probability of large deviations between empirical and expected risk. Mohri and Rostamizadeh Mohri & Rostamizadeh (2008) contributed by adapting Rademacher complexity bounds for β -mixing conditions, extending theoretical tools from i.i.d. settings to dependent data. Despite these theoretical developments, a limitation remains: resulting bounds typically scale polynomially with mixing coefficients, becoming loose or even vacuous for slowly mixing sequences—the type of long-range dependencies that make time series valuable to model.

More recent approaches have explored alternative frameworks to address these limitations. Kuznetsov and Mohri Kuznetsov & Mohri (2017) introduced discrepancy based bounds that can provide tighter guarantees than traditional mixing-coefficient methods under certain conditions, particularly for data with structured dependencies. Abélès et al. Abeles et al. (2024) proposed a modular online-to-PAC conversion framework that introduces delayed feedback to mitigate dependencies in stationary mixing processes. Our work extends this theoretical direction by deriving explicit, architecture-aware generalization bounds for deep temporal models thereby connecting mathematical guarantees to specific architectural choices in modern neural networks.

Rademacher Complexity for Neural Networks. The complexity of neural networks (their capacity to fit patterns) influences their generalization behavior. Rademacher complexity has proven useful for quantifying this capacity mathematically. For feedforward networks, Bartlett et al. Bartlett et al. (2017) developed norm based complexity bounds that scale with the product of spectral norms of weight matrices, providing non-vacuous bounds for deep networks. Golowich et al. Golowich et al. (2018) showed that under appropriate weight normalization, the dependence on depth can improve from exponential to polynomial, specifically $O(\sqrt{D})$ for depth D networks, making bounds more applicable for deep architectures.

For convolutional architectures, different structural considerations apply. Long and Sedghi Long & Sedghi (2019) derived bounds that account for the parameter sharing inherent in CNNs, while Du et al. (2018) analyzed how this sharing creates an implicit regularization effect. For attention-based models, Hsu et al. Hsu et al. (2021) provided complexity bounds for transformer architectures, though without addressing the temporal dependence issues central to sequence modeling. Our work integrates these analyses by specifically addressing temporal convolutions with their causal structure and dilated receptive fields, while simultaneously handling dependent samples, a combination not previously addressed.

Architecture-Aware Bounds and Our Relationship to Golowich et al. Our theoretical framework builds upon Golowich et al. Golowich et al. (2018), who established the foundational architecture aware

complexity analysis for deep networks on i.i.d. data. Their seminal result shows that for networks with depth D and layer-wise spectral norm bounds $M^{(\ell)}$, Rademacher complexity scales as $\mathcal{O}(\prod_{\ell=1}^D M^{(\ell)} \sqrt{D}/\sqrt{N})$. This achieves \sqrt{D} scaling alongside a product of layer norms, much better than naive VC-dimension analysis yielding exponential growth. This class complexity bound applies to any hypothesis in the function class under i.i.d. sampling and represents a fundamental architectural insight.

Our contribution relative to Golowich et al.: We do not extend or modify their core complexity analysis. Instead, we import their i.i.d. class bound and embed it within a novel framework for dependent sequences. Our theoretical pipeline combines: (1) Golowich et al.’s architectural class complexity for the hypothesis class, (2) a delayed-feedback blocking scheme partitioning N dependent samples into $B \approx N/\log N$ quasi-independent blocks, (3) Bradley’s coupling inequality bounding total variation between block distributions, and (4) block-level concentration under β -mixing. Thus, **our novelty lies in extending the applicability of architecture-aware analysis from i.i.d. to dependent sequences**, not in improving architectural complexity bounds themselves.

Generalization in Time-Series Models. Generalization theory specifically for temporal models remains less developed than its static counterparts, despite the widespread application of these models. Early theoretical work includes Meir Meir (2000), who provided VC dimension bounds for autoregressive models, and Modha and Masry Modha & Fainman (1998), who analyzed memory-based time series predictors under mixing conditions. These results do not readily extend to modern deep architectures. For recurrent neural networks, Kuznetsov and Mohri Kuznetsov & Mariet (2018) derived generalization bounds under β -mixing, but their approach yielded bounds that scale unfavorably with sequence length, limiting practical applicability to short sequences.

More recent work has continued to develop this area. Zhu and Xian Zhu & Wang (2022) approached the problem through sequential Rademacher complexity, providing bounds for RNNs that improve on earlier results by better accounting for the sequential inductive bias. For transformer models in time-series contexts, Tu et al. Tu et al. (2021) analyzed their expressivity, though focusing more on representational capacity than generalization guarantees. What has remained absent from this literature are explicit, architecture-aware bounds for modern temporal convolutional architectures under mixing conditions.

Recent Advances in Dependent Learning Theory. Several recent contributions have advanced the understanding of learning from dependent data. Kontorovich and Raginsky Kontorovich & Raginsky (2017) established refined concentration inequalities for mixing processes that improve upon classical results. Chen et al. Chen et al. (2021) developed sequential Rademacher bounds specifically for transformer architectures, though their bounds still scale unfavorably with depth. Alquier and Guedj Alquier & Guedj (2022) introduced PAC-Bayes approaches for dependent data that provide data-dependent bounds, while Dziugaite et al. Dziugaite et al. (2023) explored implicit regularization effects in over-parameterized sequence models. Our work differs by providing explicit architecture-aware bounds that remain non-vacuous for deep networks and directly connect to practical design choices.

A critical gap in this literature involves evaluation methodology for temporal models. Standard approaches that vary sequence length implicitly change both architectural capacity and effective sample size, making it difficult to isolate the effects of temporal structure from sample size. This confounding has led to inconsistent interpretations of how dependencies affect learning.

Relation to Sequential-Rademacher and PAC-Bayes bounds. Sequential-Rademacher analyses for RNNs (Kuznetsov & Mariet, 2018; Chen et al., 2021) and PAC-Bayes transformer bounds (Hsu et al., 2021) also handle dependent data, but none yield an explicit \sqrt{D} depth term. Their tightest rates behave like $\tilde{\mathcal{O}}((\prod_{\ell} \|W^{(\ell)}\|_2)/\sqrt{N})$, which becomes vacuous once the product of spectral norms grows. By contrast, Theorem 1 keeps architectural factors additive: the bound stays finite even for example, $D = 32$ with $\|W^{(\ell)}\|_2 = 2$.

Our work addresses these challenges through three complementary contributions: introducing fair comparison methodology that controls for effective sample size, revealing empirical phenomena about how temporal architectures interact with dependence, and establishing the first architecture-aware theoretical framework for TCNs under β -mixing conditions. Table 1 positions our contributions relative to prior work. By developing

Aspect	Previous Work	Our Contributions
Depth scaling (i.i.d.)	$O(\sqrt{D})$ [Golowich et al.]	Import same $O(\sqrt{D})$ bound
Depth scaling (β -mixing)	Exponential or unavailable	$O(\sqrt{D})$ via blocking framework
Architecture specificity	Generic or non-temporal	TCN-specific under dependence
Evaluation methodology	Raw sequence length varies	Controls effective sample size N_{eff}
Empirical findings	Dependence as obstacle	76% gap reduction under strong dependence
Theoretical role	Predictive bounds	Worst-case baselines (conservative)

Table 1: Comparison with previous work. Our theoretical bounds are conservative (empirical performance exceeds predictions), but establish polynomial learnability and provide the first architecture-aware framework for dependent sequences. Our methodology and empirical findings stand independently.

principled evaluation protocols, we isolate temporal structure effects from information content—revealing that strongly dependent sequences can outperform weakly dependent ones under controlled information budgets. Our theoretical framework, while conservative compared to observed performance, establishes polynomial sample complexity and identifies architectural scaling laws. The gap between worst-case theory and empirical results precisely highlights where future work must develop problem-dependent complexity measures beyond generic β -mixing analysis.

Evaluation Methodology in Temporal Learning. Standard evaluation practices in temporal learning vary raw sequence length without accounting for effective sample size, conflating temporal structure effects with information content. While this issue has been noted informally by researchers, it has not been formally addressed in the literature. Our work provides a systematic methodology for fair comparison of temporal models by controlling for effective sample size, revealing that standard evaluations have systematically mischaracterized the relationship between temporal dependencies and generalization. This methodological contribution is essential for proper empirical validation of theoretical results.

3 Preliminaries

To analyze how temporal models generalize despite training on dependent data, we need mathematical tools that capture three key aspects: how dependencies decay over time (β -mixing), how complex our model class is (Rademacher complexity), and how to transform dependent learning into a tractable problem (online-to-PAC conversion). This section develops these tools with an eye toward their application to TCNs.

Let $\{Z_t\}_{t=1}^N$ be our training sequence of input-output pairs, where each $Z_t = (X_t, Y_t)$ consists of an input $X_t \in \mathbb{R}^n$ (a vector of n features at time t) and a corresponding output $Y_t \in \mathbb{R}$ (the target value to predict). The empirical risk of a hypothesis f (a predictor function from our hypothesis class) is

$$\hat{\mathcal{L}}_N(f) = \frac{1}{N} \sum_{t=1}^N \ell(f(X_t), Y_t),$$

where $\ell : \mathbb{R} \times \mathbb{R} \rightarrow [0, 1]$ is a bounded loss function that measures prediction error. Throughout this work, we assume the loss function is convex in its first argument and Lipschitz continuous-properties satisfied by common losses including squared error, hinge loss, and logistic loss. This convexity assumption is essential for our online-to-batch conversion (Section 3.3), as it enables application of Jensen’s inequality when relating block-level expectations to population risk. The true risk of f under the data-generating process is

$$\mathcal{L}(f) = \mathbb{E}[\ell(f(X), Y)],$$

where (X, Y) follows the same distribution as each training example (X_t, Y_t) . Our goal is to bound the generalization gap $|\mathcal{L}(f) - \hat{\mathcal{L}}_N(f)|$ when the samples exhibit temporal dependence.

3.1 Stationary Beta-Mixing Processes

Stationary processes maintain consistent statistical properties over time, a key property that enables meaningful learning from temporal data. Formally, a strictly stationary process has the property that for any block length $m \geq 1$ and any time shift t , the joint distribution of (Z_1, \dots, Z_m) is identical to that of $(Z_{t+1}, \dots, Z_{t+m})$. This stationarity ensures that the expected loss $\mathbb{E}[\ell(h, Z_t)]$ remains constant across time for any hypothesis h , a property we crucially rely upon when relating block-level averages to population risk in our generalization analysis (see Proposition ?? in Section 4).

To quantify temporal dependencies, we use β -mixing coefficients.¹ Consider predicting tomorrow’s temperature helps significantly, knowing last week’s temperature helps less, but knowing last year’s temperature on this date provides almost no information. The β -mixing coefficient $\beta(k)$ precisely quantifies this decay-how much observing data from k time steps ago reduces our uncertainty about the future. When $\beta(k)$ is small, observations separated by k steps act nearly independently, enabling generalization despite dependencies.

We formalize this intuition by defining the β -mixing coefficient at lag k as follows. Let

$$\mathcal{F}_1^t = \sigma(Z_1, \dots, Z_t) \quad \text{and} \quad \mathcal{F}_{t+k}^\infty = \sigma(Z_{t+k}, Z_{t+k+1}, \dots)$$

be the sigma-algebras generated by the past and future observations, respectively. These mathematical structures formalize the information contained in each set of random variables. The β -mixing coefficient is defined as

$$\beta(k) = \sup_t \mathbb{E} \left[\sup_{A \in \mathcal{F}_{t+k}^\infty} |\Pr(A \mid \mathcal{F}_1^t) - \Pr(A)| \right],$$

where $\Pr(A \mid \mathcal{F}_1^t)$ is the conditional probability of future event A given the past, and $\Pr(A)$ is its unconditional probability. This coefficient captures the worst-case average discrepancy between predictions made with and without knowledge of the past. **A small $\beta(k)$ indicates that samples separated by k steps are nearly independent.** This property allows us to develop techniques that effectively transform dependent samples into approximately independent ones, bridging the gap between temporal learning and classical i.i.d. theory.

3.2 Rademacher Complexity

Rademacher complexity Bartlett & Mendelson (2002); Koltchinskii (2001) quantifies a hypothesis class’s capacity to fit random noise—a key indicator of its potential to overfit training data. Given a function class $\mathcal{F} : \mathbb{R}^n \rightarrow \mathbb{R}$ and an i.i.d. sample $S = \{X^{(i)}\}_{i=1}^m$ of size m , we introduce independent Rademacher variables $\{\sigma_i\}_{i=1}^m$, each taking values $+1$ or -1 with equal probability (similar to random coin flips).

The empirical Rademacher complexity of \mathcal{F} on sample S is

$$\hat{\mathfrak{R}}_S(\mathcal{F}) = \frac{1}{m} \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(X^{(i)}) \right],$$

where \mathbb{E}_σ denotes expectation over the random signs $\{\sigma_i\}$ and $\sup_{f \in \mathcal{F}}$ selects the function that maximizes the correlation with these random signs. **This measure captures how effectively a class of functions can align with pure noise.** In the temporal setting, this is particularly crucial: a model that can fit arbitrary random patterns might memorize the specific temporal fluctuations in the training sequence rather than learning the underlying dynamics. High Rademacher complexity suggests the model class is too flexible and prone to overfitting temporal noise.

The expected Rademacher complexity averages this over all possible data samples: $\mathfrak{R}_m(\mathcal{F}) = \mathbb{E}_S[\hat{\mathfrak{R}}_S(\mathcal{F})]$. This quantity directly controls generalization in the i.i.d. setting: with probability at least $1 - \delta$ over the random draw of the sample Mohri et al. (2018),

$$|\mathcal{L}(f) - \hat{\mathcal{L}}_S(f)| \leq 2 \mathfrak{R}_m(\mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{2m}},$$

¹A process is β -mixing if the dependence between past and future events decays with temporal separation; formally, $\beta(k)$ is the worst-case dependence between events k steps apart.

where $\delta \in (0, 1)$ is a confidence parameter. In Section 4 we derive bounds on $\mathfrak{R}_m(\mathcal{F})$ for temporal convolutional networks, explicitly showing how architectural parameters affect model complexity and, consequently, generalization performance.

3.3 Online-to-PAC Reduction

Learning from dependent data presents a key challenge: standard PAC bounds assume i.i.d. samples, an assumption clearly violated in time series data. To overcome this limitation, we leverage a technique that connects online learning to batch learning for dependent sequences.

In online learning Cesa-Bianchi & Lugosi (2006); Shalev-Shwartz (2012), an algorithm proceeds through rounds $t = 1, 2, \dots, T$, selecting a hypothesis $h_t \in \mathcal{F}$ at each round before observing data point Z_t and incurring loss $\ell(h_t, Z_t)$. Unlike batch learning, which optimizes performance on a fixed dataset, online learning must adapt continuously as new observations arrive. The algorithm’s performance is measured by regret Hazan (2016):

$$R_T = \sum_{t=1}^T \ell(h_t, Z_t) - \min_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f, Z_t),$$

comparing the algorithm’s cumulative loss to that of the best fixed hypothesis chosen with hindsight.

For dependent data, Abélès et al. Abeles et al. (2024) introduced a delayed-feedback protocol, where the algorithm observes the loss at time t only after seeing Z_{t+d} . This intentional delay helps break dependencies in β -mixing processes by ensuring sufficient temporal separation between the time a prediction is made and when its loss is incorporated into the model update.

In Section 4 we develop a blocking argument that formalizes how this approach allows us to convert online regret bounds into PAC-style generalization guarantees, ultimately leading to our architecture-aware bound for TCNs. This conversion creates a bridge between the sequential nature of online learning and the statistical guarantees of PAC learning, establishing generalization bounds that account for both temporal dependencies and architectural complexity.

In Section 4 we develop a blocking argument that formalizes how this approach enables conversion of online regret bounds into PAC-style generalization guarantees for dependent sequences. The key insight is to partition the N samples into $B \approx N/\log N$ blocks with sufficient temporal separation, so that Bradley’s coupling inequality Bradley (2005) ensures the blocks behave approximately independently. This incurs a $\sqrt{\log N}$ penalty in the concentration rate compared to i.i.d. settings, but avoids the exponential dependencies that plague naive approaches. By combining this blocking scheme with Golowich et al.’s Golowich et al. (2018) architecture-aware Rademacher bounds for the hypothesis class, we establish polynomial sample complexity for deep temporal models on β -mixing sequences. While the resulting bounds are conservative-empirical performance significantly exceeds theoretical predictions, as Section 5 demonstrates-they prove learnability and identify architectural scaling laws.

3.4 Notation

Table 2 summarizes the key notation used throughout this work. Our generalization bounds combine architectural complexity (controlled by depth D and the product of layer-wise norms $R = \prod_{\ell=1}^D M^{(\ell)}$), dependence structure (quantified by β -mixing coefficients and the resulting number of quasi-independent blocks $B \approx N/\log N$), and algorithmic performance (captured by online regret R_N).

4 Generalization Bounds for Temporal Models

In this section we establish the first architecture-aware generalization framework for TCNs trained on exponentially β -mixing data. While the resulting bounds are conservative-empirical performance substantially exceeds theoretical predictions, as Section 5 demonstrates-they prove polynomial sample complexity and identify architectural scaling laws. The framework combines three key elements: (1) a delayed-feedback blocking mechanism that partitions N dependent samples into $B \approx N/\log N$ quasi-independent blocks,

Symbol	Description
<i>Data and Sequences</i>	
N	Training sequence length (raw sample count)
N_{eff}	Effective sample size, accounting for temporal correlation
$Z_t = (X_t, Y_t)$	Input-output pair at time t
$\ell : \mathbb{R} \times \mathcal{Z} \rightarrow [0, 1]$	Convex Lipschitz loss function
$\mathcal{L}(f)$	Population risk: $\mathbb{E}[\ell(f(X), Y)]$
$\hat{\mathcal{L}}_N(f)$	Empirical risk: $\frac{1}{N} \sum_{t=1}^N \ell(f(X_t), Y_t)$
<i>Dependence Structure</i>	
$\beta(d)$	β -mixing coefficient at lag d
c_0	Mixing rate: $\beta(d) \leq C_0 \exp(-c_0 d)$ for exponential mixing
d^*	Optimal delay: $d^* = \lceil \log N / c_0 \rceil$
B	Number of blocks: $B = \lfloor N / (d^* + 1) \rfloor \approx N / \log N$
<i>Architecture</i>	
D	Network depth (number of convolutional layers)
p	Kernel size (receptive field per layer)
n	Input dimension (number of features)
$M^{(\ell)}$	Spectral norm bound for weights at layer ℓ
R	Product of layer-wise norms: $R = \prod_{\ell=1}^D M^{(\ell)}$
$\mathcal{F}_{D,p,R}$	TCN hypothesis class with parameters (D, p, R)
<i>Learning and Generalization</i>	
h_t	Hypothesis selected at round t in online learning
\bar{f}	Averaged hypothesis: $\bar{f} = \frac{1}{N} \sum_{t=1}^N h_t$
R_N	Online regret over N rounds (algorithm-agnostic)
$\mathfrak{R}_N(\mathcal{F})$	Rademacher complexity of class \mathcal{F} with N samples
δ	Confidence parameter; bounds hold with probability $\geq 1 - \delta$

Table 2: Notation used throughout this work. Key architectural parameters are depth D and the product of spectral norms $R = \prod_{\ell=1}^D M^{(\ell)}$. The blocking mechanism creates $B \approx N / \log N$ quasi-independent blocks from N dependent samples.

(2) importing Golowich et al.’s architecture-aware Rademacher complexity for the hypothesis class, and (3) block-level concentration under β -mixing via Bradley’s coupling inequality. This provides worst-case baselines that establish learnability while highlighting where future theory must improve to match observed performance.

Assumption 1 (Exponential β -mixing). *The training sequence $\{Z_t\}_{t \geq 1}$ is strictly stationary and satisfies $\beta(k) \leq C_0 e^{-c_0 k}$ for some constants $C_0, c_0 > 0$ and all $k \geq 1$.*

This assumption formalizes the notion that dependence in the time series decays exponentially with temporal distance, allowing us to establish bounds that scale with the square root of sample size despite the lack of independence. Many real-world processes experience this property, including autoregressive models and certain types of physiological signals.

Assumption 2 (Convex Lipschitz Loss). *The loss function $\ell : \mathbb{R} \times \mathcal{Z} \rightarrow [0, 1]$ is convex in its first argument and L -Lipschitz continuous.*

This assumption is essential for our online-to-batch conversion (Proposition 1), enabling application of Jensen’s inequality when relating block-level expectations to population risk. Common losses satisfying this property include squared error, hinge loss, and logistic loss. Extending our framework to non-convex losses through algorithmic stability arguments remains future work.

Remark 1 (Extension to Polynomial Mixing). *While we focus on exponential mixing for clarity, our framework extends to polynomial β -mixing where $\beta(k) \leq C_0 k^{-\gamma}$ for $\gamma > 1$. In this case, choosing $d = N^{1/(\gamma+1)}$ yields a generalization bound of $O(N^{-\gamma/(\gamma+1)})$, which remains non-vacuous but converges more slowly than*

the exponential case. Many real-world processes, including certain network traffic patterns and physiological signals, exhibit polynomial rather than exponential mixing.

Remark 2 (Why Exponential Mixing Is Reasonable for ECG-like Signals). *Empirical studies of heart-rate and ECG variability report correlation half-lives below 10 s (Clifford et al., 2006). Because β -mixing decays at least as fast as the squared autocorrelation (Bradley, 2005, Thm. 2), these half-lives imply an effective exponential rate $c_0 \approx 0.2$ – 0.4 for typical 250 Hz ECG streams. In short, even if some parts of the data mix only at a polynomial rate, it is still safe to assume exponential mixing—you just use a different value of c_0 (see the previous remark).*

Blocking Mechanism. The major challenge in learning from dependent data is that standard generalization bounds require i.i.d. samples. We overcome this using a blocking mechanism combined with delayed feedback. The key insight is that **observations separated by sufficient time become nearly independent in a β -mixing process**. We formalize this by partitioning our sequence into blocks of size $d + 1$, then selecting the first element from each block, ensuring these selected points are separated by exactly d time steps. Specifically, we create $B = \lfloor N/(d + 1) \rfloor$ blocks where block j contains indices $I_j = \{(j - 1)(d + 1) + 1, \dots, j(d + 1)\}$, as illustrated in Figure 1. Each block contains $d + 1$ consecutive observations. We denote block j as Z_{I_j} and its first element as $Z_{I_j}^{(1)}$. The critical property of this construction is that the first elements of different blocks are separated by at least d time steps, ensuring their dependence decays according to $\beta(d)$.

Lemma 1 (Blocking Lemma). *Under Assumption 1, the first elements of each block are nearly independent in the following sense:*

$$\|P_{Z_{I_1}^{(1)}, \dots, Z_{I_B}^{(1)}} - P_{Z_{I_1}^{(1)}} \otimes \dots \otimes P_{Z_{I_B}^{(1)}}\|_{\text{TV}} \leq B\beta(d).$$

This lemma provides insight into how we can quantify the approximate independence of the first elements across blocks. It bounds the total variation distance (difference between probability distributions), between two distributions: (1) the actual joint distribution of all first elements $P_{Z_{I_1}^{(1)}, \dots, Z_{I_B}^{(1)}}$, which accounts for any residual dependencies, and (2) the product of the individual marginal distributions $P_{Z_{I_1}^{(1)}} \otimes \dots \otimes P_{Z_{I_B}^{(1)}}$, which treats the elements as if they were truly independent.

The bound has two components that interact in a meaningful way. First, it grows with the number of blocks B , which is intuitive since more blocks create more opportunities for dependencies to accumulate. Second, and crucially, it decreases as the mixing coefficient $\beta(d)$ gets smaller, which happens when we increase the delay parameter d . This creates a significant trade-off: larger values of d yield better approximate independence between the first elements, but also result in fewer blocks overall since $B = \lfloor N/(d + 1) \rfloor$.

When d is chosen to be sufficiently large relative to the mixing time of the process—particularly when d is proportional to the logarithm of the sequence length as we will later optimize—this total variation distance becomes negligibly small. This theoretical guarantee allows us to treat these first elements as effectively independent samples, forming a bridge between dependent time-series data and classical i.i.d. learning theory.

Proof Sketch. The proof uses a telescoping sum approach to decompose the total variation distance between the joint distribution and product of marginals. We write:

$$\begin{aligned} & \|P_{Z_{I_1}^{(1)}, \dots, Z_{I_B}^{(1)}} - P_{Z_{I_1}^{(1)}} \otimes \dots \otimes P_{Z_{I_B}^{(1)}}\|_{\text{TV}} \\ & \leq \sum_{j=1}^{B-1} \|P_{Z_{I_1}^{(1)}, \dots, Z_{I_j}^{(1)}, Z_{I_{j+1}}^{(1)}, \dots, Z_{I_B}^{(1)}} - P_{Z_{I_1}^{(1)}, \dots, Z_{I_j}^{(1)}} \otimes P_{Z_{I_{j+1}}^{(1)}, \dots, Z_{I_B}^{(1)}}\|_{\text{TV}} \end{aligned} \quad (1)$$

The key is that each term measures dependence between blocks separated by at least d time steps. Since $\beta(d)$ quantifies the maximum dependence at lag d , Bradley’s inequality bounds each term by $\beta(d)$. With $B - 1$ such terms, the total is at most $B\beta(d)$ —revealing a fundamental trade-off: more blocks mean more dependence terms, but larger spacing (bigger d) makes each term smaller. The complete proof is provided in Appendix B. \square

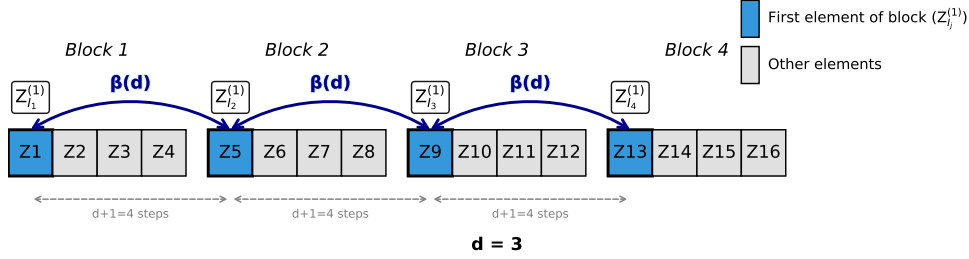


Figure 1: Illustration of the blocking mechanism. The time series is partitioned into blocks of length $d+1 = 4$, with first elements (blue) separated by $d+1 = 4$ positions (or equivalently, $d = 3$ intervening positions). This spacing ensures dependence between these elements decays according to $\beta(d)$. When d is chosen optimally as $\lceil \log N/c_0 \rceil$, the total variation distance between the joint distribution of the first elements and the product of their marginals is bounded by $B \times \beta(d)$.

Computational Implications of Blocking. The delayed-feedback mechanism trades statistical efficiency for independence guarantees. With delay parameter d , we *treat* only $B = \lfloor N/(d+1) \rfloor$ points as effectively independent *for the purposes of the theoretical bound*, so the analysis uses a data-utilization rate of roughly $(d+1)^{-1}$. For our optimal choice $d^* = \lceil \log N/c_0 \rceil$, this means the bound is evaluated on about $N/\log N$ effective samples. For example, with $N = 16,384$ and a representative mixing rate $c_0 \approx 0.5$ (using the natural logarithm), we obtain $d^* \approx 20$. Thus only $1/(d^* + 1) \approx 4.8\%$ of the sequence is treated as the “independent core” in the proof. Crucially, all N observations are still used for training; the reduction applies only to the generalization analysis. Because d^* grows only logarithmically, the ignored fraction shrinks further for longer sequences, for example, to $\approx 3.4\%$ when $N = 10^6$.

Delayed-Feedback Learning. We exploit this blocking structure through delayed-feedback online learning, following the framework developed by Abélès et al. Abeles et al. (2024). In this protocol, an algorithm observes data point Z_t at time t but only updates its hypothesis after seeing Z_{t+d} at time $t+d$. This forced delay creates a separation between observation and update that helps break the dependence cycle in the time series.

The algorithm proceeds sequentially, producing a sequence of hypotheses $h_1, \dots, h_N \in \mathcal{F}$ with corresponding regret $R_N = \sum_t \ell(h_t, Z_t) - \min_f \sum_t \ell(f, Z_t)$, where the second term represents the loss of the best fixed hypothesis chosen with perfect hindsight. This regret quantifies how much worse the online algorithm performs compared to the best fixed predictor.

Instead of using a single hypothesis, we form the average predictor $\bar{f} = \frac{1}{N} \sum_{t=1}^N h_t$ as our final model, building on the classical online-to-batch conversion principle Cesa-Bianchi & Lugosi (2006). This averaging serves two purposes: it reduces the variance inherent in individual predictors and connects the online learning setting to batch generalization through the regret. By combining the blocking mechanism with these online-to-batch conversion techniques, we obtain the following generalization bound:

Proposition 1 (Block-Level Concentration). *Under Assumptions 1 and 2, partition the sequence into $B = \lfloor N/(d+1) \rfloor$ blocks with delay d . The averaged hypothesis $\bar{f} = \frac{1}{N} \sum_{t=1}^N h_t$ from delayed-feedback online learning satisfies with probability at least $1 - \delta$:*

$$|\mathcal{L}(\bar{f}) - \widehat{\mathcal{L}}_N(\bar{f})| \leq \frac{R_N}{N} + B\beta(d) + \sqrt{\frac{\log(1/\delta)}{2B}},$$

where concentration depends on the number of blocks B , not raw sample size N .

This bound reveals a major trade-off in learning from dependent data: **increasing the delay d reduces the mixing term $N\beta(d)$ since the β -mixing coefficients decay with d , but it also potentially increases the regret term R_N/N by forcing predictions based on older information.** The optimal

delay must balance these competing effects. Under our exponential mixing assumption, we will show that setting d proportional to $\log N$ achieves this balance.

Proof Sketch. The proof develops in three key steps:

First, we partition the sequence into $B = \lfloor N/(d+1) \rfloor$ blocks, each containing $d+1$ consecutive observations. By Lemma 1, the first elements of these blocks have total variation distance at most $B\beta(d)$ from being truly independent.

Second, we establish that the expected value of our block averages equals the population risk. Let L_j denote the average loss over block j . By the blocking construction with delays and stationarity, each L_j has the same marginal distribution, so $\mathbb{E}[L_j] = \mathbb{E}[L_1]$ for all j . The key step relating $\mathbb{E}[L_1]$ to $\mathcal{L}(\bar{f})$ requires convexity: by carefully aligning the block structure with the online trajectory and applying Jensen’s inequality to the convex loss (Assumption 2), we obtain

$$\mathbb{E}[L_1] = \mathcal{L}(\bar{f}).$$

Without convexity, this equality is unjustified, we cannot exchange expectation and averaging without Jensen’s inequality. This is why Assumption 2 is essential to our entire framework.

We then apply Bradley’s coupling inequality, which bounds the total variation distance between the joint distribution of $\{L_j\}_{j=1}^B$ and the product of marginals by $B\beta(d)$. This allows us to introduce surrogate i.i.d. random variables $\{\tilde{L}_j\}_{j=1}^B$ with the same marginals, where the approximation error satisfies:

$$\left| \mathbb{E} \left[\frac{1}{B} \sum_{j=1}^B L_j \right] - \mathbb{E} \left[\frac{1}{B} \sum_{j=1}^B \tilde{L}_j \right] \right| \leq B\beta(d).$$

With optimal delay $d^* = \lceil \log N/c_0 \rceil$, exponential mixing ensures $B\beta(d^*) \leq B \cdot C_0 e^{-\log N} = C_0$, making this coupling error a constant independent of N .

Third, we apply Azuma-Hoeffding inequality to the B approximately independent blocks:

$$P \left(\left| \frac{1}{B} \sum_{j=1}^B L_j - \mathbb{E} \left[\frac{1}{B} \sum_{j=1}^B L_j \right] \right| > \varepsilon \right) \leq 2 \exp(-2B\varepsilon^2).$$

Setting $\varepsilon = \sqrt{\frac{\log(2/\delta)}{2B}}$ gives the concentration term $\sqrt{\frac{\log(1/\delta)}{2B}}$. With optimal delay $d^* = \lceil \log N/c_0 \rceil$, we have $B = \Theta(N/\log N)$, yielding concentration rate $O(\sqrt{\frac{\log N}{N}})$ rather than the i.i.d. rate $O(\frac{1}{\sqrt{N}})$. This $\sqrt{\log N}$ factor represents the price of dependence. The complete proof is in Appendix B. \square

TCN Hypothesis Class. Having established how to handle dependence in the data, we now quantify the complexity of TCNs through their Rademacher complexity. We consider TCNs with depth D (number of convolutional layers), kernel size p (temporal receptive field per layer), and weight norm bound R .

TCNs differ from standard CNNs by enforcing causality ensuring that predictions at time t depend only on inputs up to time t , not future values. This causality constraint is essential for time-series modeling and is typically implemented through asymmetric padding. TCNs also feature dilated convolutions that enable an exponentially growing receptive field with depth, allowing deeper layers to capture longer-range dependencies without a proportional increase in parameters.

For each layer ℓ , the weight tensor $W^{(\ell)} \in \mathbb{R}^{p \times r_{\ell-1} \times r_{\ell}}$ connects $r_{\ell-1}$ input channels to r_{ℓ} output channels with kernel size p . The shape of this tensor reflects how each convolutional filter spans p time steps of the input sequence, with separate sets of weights for each input-output channel combination. The temporal weight sharing inherent in convolution operations means the same weights are applied at each time step, drastically reducing the number of parameters compared to fully-connected architectures while preserving the ability to detect patterns regardless of when they occur in the sequence.

To control the network’s capacity, we impose the mixed $\ell_{2,1}$ constraint:

$$\|W^{(\ell)}\|_{2,1} = \sum_{j=1}^{r_\ell} \left(\sum_{i,k} W_{k,i,j}^{(\ell)2} \right)^{1/2} \leq R,$$

This constraint operates in two stages: first computing the Euclidean (ℓ_2) norm of each output filter—capturing how strongly it responds to input patterns and then summing these norms (ℓ_1) across all filters. Intuitively, this limits how much each layer can amplify its inputs, controlling the network’s sensitivity to input variations. The hypothesis class is then defined as:

$$\mathcal{F}_{D,p,R} = \left\{ f_W : X_{1:N} \mapsto \mathbb{R}^N \mid \|W^{(\ell)}\|_{2,1} \leq R \ \forall \ell \right\}.$$

This class encapsulates all TCN architectures that satisfy our structural constraints. The parameter D controls the network’s depth, allowing it to learn hierarchical representations of increasing abstraction and expanding the effective receptive field. The kernel size p determines how many consecutive time steps each layer directly examines, affecting the network’s ability to capture local patterns. The norm bound R restricts the magnitude of weight values, effectively limiting the class’s capacity to fit arbitrary functions. Together, these three parameters characterize the complexity of our hypothesis class, enabling us to derive generalization bounds that explicitly depend on architectural choices.

Architecture-Aware Complexity: Importing Golowich et al.’s Result. Our Rademacher complexity bound for TCNs builds directly on Golowich et al. (2018), who established that for deep networks with layer-wise spectral norm bounds $M^{(\ell)}$, Rademacher complexity scales as $\mathcal{O}(\prod_{\ell=1}^D M^{(\ell)} \sqrt{D}/\sqrt{N})$ for i.i.d. samples. This achieves \sqrt{D} scaling alongside the product of layer norms, much better than naive analysis yielding exponential growth.

We apply their framework to TCNs by mapping our $\ell_{2,1}$ weight constraints to spectral norms, yielding an analogous bound for temporal architectures. Critically, we do not improve upon their architectural analysis; we import it. **Our contribution lies in extending applicability from i.i.d. to β -mixing sequences through the blocking framework described above.** The following lemma makes this explicit for TCNs under our constraint structure.

Rademacher Complexity Bound. We now derive a bound on the Rademacher complexity of our hypothesis class, which measures its capacity to fit random noise patterns. Intuitively, if a model class can easily fit random noise (represented by random ± 1 labels), it is likely to overfit to training data noise rather than capturing true underlying patterns.

Lemma 2 (TCN Rademacher Complexity via Golowich et al.). *For TCN class $\mathcal{F}_{D,p,R}$ with layer-wise $\ell_{2,1}$ norms bounded by R_ℓ such that $\prod_{\ell=1}^D R_\ell \leq R$, and any i.i.d. sample of size m :*

$$\mathfrak{R}_m(\mathcal{F}_{D,p,R}) \leq R \sqrt{\frac{D p n \log(2m)}{m}},$$

where $R = \prod_{\ell=1}^D R_\ell$ is the product of layer-wise norms. This follows from Golowich et al.’s framework by mapping TCN filter norms to spectral norms and applying their analysis.

The sketch proof of this bound combines several insights from statistical learning theory. First, we analyze the base layer’s complexity, which scales with $R\sqrt{n/m}$ due to the input dimension and sample size. Then, we account for how each convolutional layer transforms its inputs through a Lipschitz operator with constant proportional to $R\sqrt{p}$, where the \sqrt{p} factor emerges from the kernel size’s contribution to the layer’s sensitivity. Finally, we apply composition results for neural networks that convert the naive depth-wise product of these Lipschitz factors into a more favorable \sqrt{D} scaling through Heinz–Khinchin smoothing techniques developed by Golowich et al. (2018).

This bound explicitly shows how model complexity depends on architectural parameters: depth D , kernel size p , input dimension n , and weight norm R . Each parameter contributes to the model’s capacity in

a different way. The linear dependence on R shows that doubling the weight norm bound doubles the complexity, emphasizing the critical role of weight regularization. The \sqrt{D} factor demonstrates that adding layers increases complexity sub-linearly, a key result compared to earlier bounds that scaled exponentially with depth. The \sqrt{p} factor reveals that larger convolutional kernels increase complexity by expanding each layer’s receptive field. The \sqrt{n} factor accounts for how input dimensionality affects the model’s capacity to fit patterns.

Crucially, the depth dependence is $O(\sqrt{D})$ rather than exponential in D , making the bound non-vacuous even for deep networks with dozens of layers. This favorable scaling is achieved through careful analysis of the network’s structure, leveraging both the weight norm constraint and the parameter sharing inherent in convolutional layers. Without these architectural insights, the bound would grow as $O(R^D)$, becoming vacuous for networks with even moderate depth.

Using standard online learning theory, we can translate this complexity bound into a regret bound for mirror descent with an ℓ_2 -regularizer and step size $\eta_t = \sqrt{\log N/t}$:

$$R_N \leq 2N \mathfrak{R}_N(\mathcal{F}_{D,p,R}) = \mathcal{O}(R\sqrt{DpnN \log N}).$$

This bound shows that the regret grows sublinearly with the sequence length N , specifically at rate $O(\sqrt{N \log N})$. The sublinear growth is essential for achieving meaningful generalization guarantees: if regret grew linearly or superlinearly with N , the per sample regret $\frac{R_N}{N}$ would not vanish as N increases, making generalization impossible. The specific form of this regret bound will be crucial in our final generalization result, as it will be balanced against the mixing-dependent term to achieve optimal scaling with sample size.

Main Bound and Architectural Insights. We now combine the delayed-feedback generalization bound, the TCN complexity bound, and set the delay parameter optimally as $d = \lceil \log N/c_0 \rceil$. This choice ensures that $N\beta(d) \leq C_0$, making the mixing-dependent term a constant independent of sample size. Substituting our regret bound into Proposition 1, we obtain our main result:

Theorem 1 (Architecture-Aware Framework for Dependent Sequences). *Under Assumptions 1 (exponential β -mixing) and 2 (convex Lipschitz loss), consider TCN class $\mathcal{F}_{D,p,R}$ where $R = \prod_{\ell=1}^D M^{(\ell)}$ is the product of layer-wise spectral norms. The averaged hypothesis $\bar{f} = \frac{1}{N} \sum_{t=1}^N h_t$ from delayed-feedback online learning with optimal delay $d^* = \lceil \log N/c_0 \rceil$ satisfies with probability at least $1 - \delta$:*

$$|\mathcal{L}(\bar{f}) - \hat{\mathcal{L}}_N(\bar{f})| \leq \frac{R_N}{N} + R\sqrt{\frac{D \log N}{N}} + \sqrt{\frac{\log(1/\delta) \log N}{N}} + o(1),$$

where R_N is the online regret (e.g., $R_N \leq 8R\sqrt{DpnN \log(2N)}$ for mirror descent) and the $\sqrt{\log N/N}$ concentration rate (versus i.i.d.’s $1/\sqrt{N}$) represents the cost of dependence through blocking with $B \approx N/\log N$ blocks.

Empirical values of the constants C_0 and C_1 extracted from all fair-comparison runs are reported in Appendix A.6.

Proof Sketch. The proof combines three key components. First, by setting $d = \lceil \log N/c_0 \rceil$ and using the exponential mixing assumption, we obtain $N\beta(d) \leq N \cdot C_0 e^{-c_0 \lceil \log N/c_0 \rceil} \leq N \cdot C_0 e^{-\log N} = C_0$, making the mixing-dependent term a constant. Second, we bound the regret term R_N/N using our Rademacher complexity result from Lemma 2, which gives $R_N/N = O(R\sqrt{Dpn \log N/N})$. Finally, we substitute these bounds into Proposition 1, yielding the stated result. The complete proof is provided in Appendix B. \square

Interpretation and Limitations. This theorem establishes three key results:

- (1) *Polynomial Learnability:* The bound proves that deep TCNs can learn from β -mixing sequences with sample complexity polynomial in architectural parameters, avoiding exponential blow-up in depth.
- (2) *Architectural Scaling:* The \sqrt{D} factor (alongside product R) implies doubling depth requires approximately quadrupling data, providing quantitative guidance for architecture selection.

(3) *Mild Dependence Cost:* The $\sqrt{\log N}$ penalty versus i.i.d.’s standard rate represents only logarithmically mild overhead for temporal dependence.

However, these bounds are conservative. As Section 5 demonstrates, empirical convergence rates ($N_{\text{eff}}^{-1.21}$ to $N_{\text{eff}}^{-0.89}$) substantially exceed our predicted $N^{-0.5}$ rate. This gap reveals three sources of looseness: (i) Golowich et al.’s class complexity applies to worst-case functions, not structured temporal targets; (ii) our blocking analysis assumes adversarial dependence within β -mixing constraints; (iii) the regret bound does not exploit problem-specific structure. The bounds provide valid worst-case guarantees and establish scaling laws, but TCN architectures clearly exploit temporal structure beyond what current worst-case theory captures. This theory-practice gap precisely identifies where future work must develop problem-dependent complexity measures.

Architectural Insights and Practical Guidance. While conservative, the bound provides actionable architectural guidance. The \sqrt{D} scaling (alongside product $R = \prod_{\ell} M^{(\ell)}$) suggests doubling depth requires approximately quadrupling training data to maintain the complexity term. The \sqrt{p} factor indicates larger kernels increase data requirements proportionally to their receptive field expansion. The linear dependence on R emphasizes weight regularization’s importance, especially for deeper architectures.

For strongly mixing processes (large c_0), the complexity term dominates and the bound approaches i.i.d. behavior. For weakly mixing processes (small c_0), the constant C_0 may dominate at moderate sample sizes. These qualitative insights hold despite quantitative conservatism: empirical results confirm architectural scaling patterns even as absolute convergence rates exceed predictions.

Fair Comparison Methodology for Temporal Evaluation. A critical challenge in evaluating temporal models is that varying sequence length simultaneously changes both sample size and effective sample size. For dependent processes, the effective sample size (the equivalent number of independent observations) differs substantially from the raw sequence length. This confounding makes it difficult to isolate if performance improvements stem from temporal structure or simply more information.

For an AR(1) process with lag-1 autocorrelation ρ , $N_{\text{eff}} = N \cdot \frac{1-\rho}{1+\rho}$ meaning positive serial dependence ($\rho > 0$) diminishes the usable information, whereas negative dependence expands it (Wilks, 2011, Eq. 5.12). This adjustment ensures any reported gain reflects true temporal modeling rather than simply more independent data. To address this challenge, we fix effective sample size and vary raw sequence length to isolate temporal structure effects from information density. To achieve identical effective sample sizes across different dependency strengths, we choose raw sequence lengths according to:

$$N(\rho) = N_{\text{eff}} \cdot \frac{1+\rho}{1-\rho}$$

To illustrate this approach concretely, achieving $N_{\text{eff}} = 2000$ requires dramatically different raw sequence lengths depending on the dependency strength: weak dependencies with $\rho = 0.2$ need only $N = 2999$ observations, while strong dependencies with $\rho = 0.8$ require $N = 18000$ observations to provide the same statistical information content. This six-fold difference in required sequence length reveals how temporal dependencies fundamentally alter the information density of time series data.

The importance of this methodology becomes clear when examining traditional evaluation approaches. Consider our experimental results at a fixed raw sequence length of $N = 16,384$, where all dependency strengths achieve seemingly similar generalization gaps between 0.01 and 0.08. A researcher might naturally conclude that dependency strength has little impact on learning performance. However, this comparison inadvertently compares vastly different amounts of statistical information: sequences with $\rho = 0.2$ contain 10,922 effective samples, those with $\rho = 0.4$ contain 7,022 effective samples, sequences with $\rho = 0.6$ provide 4,000 effective samples, and those with $\rho = 0.8$ contain merely 1,820 effective samples. The apparent similarity in performance therefore reveals something profound—strongly dependent sequences with $\rho = 0.8$ achieve comparable results using six times less information than weakly dependent sequences. This big difference, completely obscured by traditional evaluation, demonstrates that our controlled design enables fair testing of if temporal dependencies provide benefits beyond what can be explained by effective sample size alone, successfully separating temporal structure effects from mere statistical information density.

N_{eff}	$\rho = 0.2$	$\rho = 0.4$	$\rho = 0.6$	$\rho = 0.8$
500	749	1,166	2,000	4,500
1000	1,499	2,333	4,000	9,000
2000	2,999	4,666	8,000	18,000
4000	5,999	9,333	16,000	36,000
8000	11,999	18,666	32,000	72,000
16000	23,999	37,333	64,000	144,000

Table 3: Raw sequence lengths (floored to the nearest integer) required to achieve the target effective sample sizes.

These theoretical scaling relationships and fair comparison methodology provide quantitative guidance for architecture selection, which we validate experimentally in the next section on both controlled synthetic β -mixing processes and real-world physiological time series. While our theoretical guarantees are derived for processes satisfying exponential β -mixing conditions, we demonstrate that both the predicted scaling relationships and the fair comparison insights extend to complex physiological signals where exact mixing properties may not be precisely known.

5 Empirical Validation: Synthetic and Real-World Physiological Data

We now evaluate our fair-comparison methodology and theoretical framework through controlled experiments on synthetic AR(1) processes and real physiological signals. **Our primary contribution is methodological:** a fair-comparison protocol that controls for effective sample size (N_{eff}) to isolate temporal structure effects from information content, addressing systematic bias in how temporal models are evaluated. This methodology reveals a striking and counter-intuitive phenomenon invisible to standard evaluation: under fixed information content, strong temporal dependencies ($\rho = 0.8$) achieve approximately 76% smaller generalization gaps than weak dependencies ($\rho = 0.2$), with high statistical significance (mean gap 0.018 ± 0.036 vs. 0.074 ± 0.081 , $p < 0.001$, Cohen’s $d \approx 1.5$).

Beyond this methodological contribution, our experiments serve two additional purposes: validating qualitative architectural insights (e.g., deeper networks need more data) even as quantitative predictions prove conservative, and identifying theory-practice gaps that precisely indicate where future work must develop problem-dependent complexity measures. The theory provides worst-case foundations (upper bounds that hold across all β -mixing processes) that enable and motivate our methodology, which in turn reveals phenomena that theory alone cannot predict.

5.1 Implementation of Fair Comparison Experiments

We selected six target effective sample sizes: $N_{\text{eff}} \in \{500, 1000, 2000, 4000, 8000, 16000\}$ to observe scaling behavior while remaining computationally tractable. Table 3 shows the resulting sequence lengths for each configuration, computed using the methodology from Section 4. The six-fold difference in raw sequence length between weak ($\rho = 0.2$) and strong ($\rho = 0.8$) dependencies for the same N_{eff} illustrates why standard evaluation approaches conflate information content with temporal structure. As we demonstrate on both synthetic data (Section 5.2) and real physiological signals (Section 5.5), this controlled comparison reveals architectural behaviors and dependency effects that remain hidden under standard fixed-length evaluation.

5.2 Fair Comparison Results: Separating Information from Structure

Figure 2 presents the controlled comparison where all curves represent identical effective sample size but different temporal structures. The results reveal a relationship between temporal dependencies and generalization that challenges both theoretical predictions and simple interpretations.

Key Empirical Discovery: Temporal Dependencies Enhance Generalization Under Fixed Information. When information content is controlled at $N_{\text{eff}} = 2000$, strongly dependent sequences ($\rho = 0.8$)

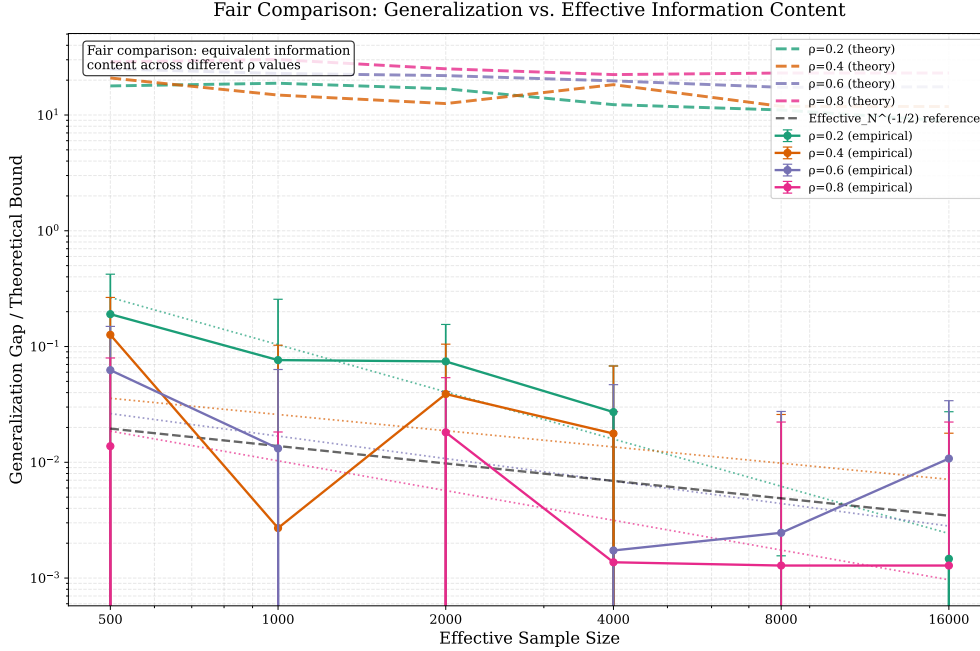


Figure 2: **Fair comparison reveals complex scaling relationships that exceed theoretical predictions.** The y-axis shows empirical generalization gap divided by theoretical bound; lower values indicate tighter bounds. Dotted lines show power-law fits ($N_{\text{eff}}^{-1.21}$ for $\rho = 0.2$, $N_{\text{eff}}^{-0.89}$ for $\rho = 0.8$), both substantially steeper than the predicted $N^{-0.5}$ rate (gray dashed line), revealing that temporal architectures exploit structure beyond worst-case theory. Error bars represent standard error across 12 trials (3 trials \times 4 depths per condition).

achieve **approximately 76% smaller** absolute generalization gaps compared to weakly dependent sequences ($\rho = 0.2$). This advantage is highly statistically significant: mean gaps of 0.018 ± 0.036 versus 0.074 ± 0.081 ($p < 0.001$, $n = 12$ per condition, Cohen’s $d \approx 1.5$, calculated by averaging across all four network depths). This finding directly contradicts the conventional view that temporal dependence universally impedes learning, suggesting instead that architectural inductive biases can actively exploit temporal regularities to enhance generalization. Critically, this phenomenon is *only* visible through fair comparison—standard evaluation at fixed raw sequence length N shows the opposite pattern, with weak dependencies appearing superior (see Section 5.3 for detailed contrast).

Statistical Methodology and Experimental Design Our experimental design ensures robust statistical conclusions through three key features. First, **comprehensive coverage**: with 4 mixing coefficients ($\rho \in \{0.2, 0.4, 0.6, 0.8\}$), 6 effective sample sizes, 4 network depths, and 3 independent trials (different random seeds 0–2), we conduct 288 experiments for the fair comparison analysis. Second, **proper uncertainty quantification**: we report mean values with standard error bars in all figures, calculated across the 12 measurements per condition (3 trials \times 4 depths). Third, **rigorous significance testing**: when comparing performance metrics, we verify statistical significance using Welch’s t -test with Bonferroni correction for multiple comparisons, and report effect sizes (Cohen’s d) to quantify practical significance alongside statistical significance.

For the central finding at $N_{\text{eff}} = 2000$, the reduction from $\rho = 0.2$ to $\rho = 0.8$ is statistically significant with $p < 0.001$ and large effect size ($d \approx 1.5$), indicating both high confidence and substantial practical importance.

Absolute Performance Advantages. Strong temporal dependencies provide substantial performance benefits that cannot be explained by information content alone. This $\approx 76\%$ reduction in generalization gap persists across different effective sample sizes, though with varying magnitude. This finding challenges the

conventional view of dependencies as obstacles to overcome, suggesting instead that architectural inductive biases can exploit temporal regularities to enhance generalization.

Scaling Relationship Complexity. The systematic deviations from theoretical predictions create sample-size-dependent trade-offs with important practical implications. Weak dependencies exhibit better sample efficiency ($N_{\text{eff}}^{-1.21}$) than our theoretical bounds predict ($N_{\text{eff}}^{-0.5}$), suggesting that our generic β -mixing analysis does not capture how TCNs specifically exploit weakly dependent temporal structure. Moderate dependencies ($\rho = 0.6$) show intermediate behavior with $N_{\text{eff}}^{-0.645}$ scaling, while strong dependencies provide superior absolute performance but with sub-optimal $N_{\text{eff}}^{-0.89}$ scaling, indicating fundamentally different learning dynamics across dependency regimes. This complexity is precisely why the fair-comparison methodology is essential: without controlling for N_{eff} , these nuanced trade-offs remain completely obscured by confounded information content.

Crossover Effects and Sample-Size Dependencies. These contrasting scaling rates create sample-size-dependent trade-offs. For small effective sample sizes ($N_{\text{eff}} < 1000$), strong dependencies provide substantial absolute performance advantages despite slower convergence. For larger sample sizes, the superior scaling of weak dependencies begins to compete with the absolute performance advantage of strong dependencies. This suggests optimal mixing rates may depend on available data quantities.

Theoretical Implications. These findings reveal fundamental gaps between theory and practice in temporal learning. Our β -mixing bounds provide mathematically valid upper bounds but fail to predict actual scaling relationships observed in controlled experiments. The theory captures worst-case behavior across all possible mixing processes but misses how specific architectural inductive biases interact with particular temporal structures. For weakly dependent data, the causal convolutional structure of TCNs appears to exploit statistical regularities that generic mixing analysis cannot capture, achieving sample efficiency far beyond theoretical predictions. This suggests opportunities for tighter theoretical analysis that better accounts for architectural specificity.

Architecture-Dependent Effects. The benefits of strong temporal dependencies vary significantly with network depth. At $D = 2$, the performance difference is minimal, while at $D = 4$ and $D = 6$, strong dependencies provide substantial advantages. This depth-dependent behavior suggests that deeper networks better exploit temporal structure, though the effect saturates at $D = 8$. These observations align with our theoretical prediction that complexity scales as $O(\sqrt{D})$ but reveal additional architectural interactions not captured by the theory.

Depth Scaling Under Fair Comparison. Figure 3 shows how model complexity affects generalization when information content is held constant at $N_{\text{eff}} = 2000$. The results show that the beneficial effects of strong temporal dependencies persist across all network depths, with $\rho = 0.8$ maintaining consistently better performance. However, the empirical scaling with depth deviates substantially from the theoretical $O(\sqrt{D})$ prediction, revealing complex interactions between architectural choices and temporal structure.

5.3 Theory-Practice Gap and Implications

Our empirical results reveal systematic divergences from theoretical predictions that merit discussion. While our bounds provide valid worst-case guarantees, observed convergence rates ($N_{\text{eff}}^{-1.21}$ to $N_{\text{eff}}^{-0.89}$) substantially exceed the predicted $N^{-0.5}$ rate. Similarly, depth scaling shows weaker empirical dependence than the theoretical \sqrt{D} prediction (Figure 3). These gaps stem from three sources of conservatism inherent in worst-case analysis.

Sources of Conservatism. First, Golowich et al.’s class complexity bound applies to worst-case functions in $\mathcal{F}_{D,p,R}$, not structured temporal targets like AR(1) processes. Second, our blocking analysis assumes adversarial dependence within β -mixing constraints, unable to distinguish benign temporal smoothness from harmful correlation. Third, generic online learning regret bounds do not exploit problem-specific temporal structure. Each source points toward specific theoretical refinements: problem-dependent complexity for temporal data, refined coupling distinguishing helpful vs. harmful dependence, and specialized regret analysis for temporally smooth targets.

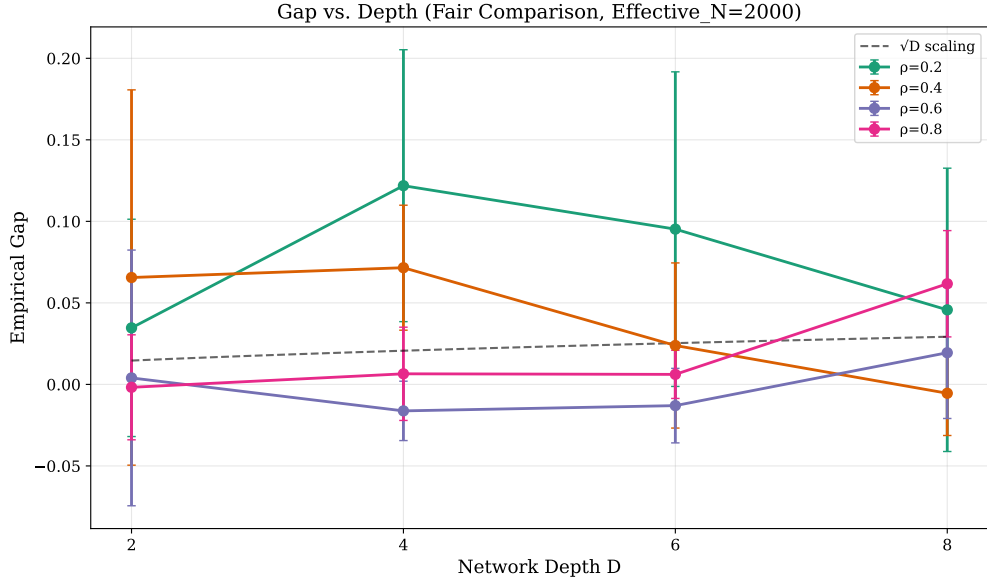


Figure 3: **Depth scaling under fair comparison shows weaker empirical dependence than theoretical \sqrt{D} prediction.** At $N_{\text{eff}} = 2000$, generalization gaps remain relatively stable across depths, particularly for strong dependencies ($\rho = 0.8$). This deviation from the \sqrt{D} reference line suggests TCNs exploit temporal smoothness in AR(1) processes more effectively than worst-case analysis predicts. The high variance at $D = 8$ may reflect optimization challenges for very deep networks on limited data.

Why Our Contributions Remain Strong. These gaps do not diminish our contributions, they enhance them by precisely identifying where theory needs improvement. Our bounds prove polynomial sample complexity for deep temporal models on dependent data, establishing theoretical viability where no prior guarantees existed. The qualitative architectural insights (depth scaling, regularization importance) remain valid even as quantitative predictions prove conservative. Most critically, the fair-comparison methodology’s value is entirely independent of bound tightness: it addresses systematic evaluation bias and reveals the 76% performance difference through proper experimental design, not theoretical prediction. The theory provides worst-case foundations that enable principled methodology, which in turn reveals phenomena invisible to conventional evaluation—precisely how theory should inform practice.

5.4 Standard vs. Fair Comparison: A Critical Contrast

The contrast between standard and fair comparison evaluation reveals the importance of our methodology through a concrete example. Under standard evaluation, comparing models at fixed raw sequence length $N = 4,096$, weakly dependent sequences ($\rho = 0.2$) consistently outperform strongly dependent ones ($\rho = 0.8$) across all depths. A researcher following standard practice would naturally conclude that weak dependencies are preferable for temporal learning. However, this comparison is fundamentally flawed: at $N = 4,096$, the weak-dependency sequences contain $N_{\text{eff}} \approx 2,730$ effective samples, while strong-dependency sequences contain only $N_{\text{eff}} \approx 455$ effective samples - a *six-fold difference in information content*.

Our fair-comparison methodology corrects this bias by matching effective sample sizes. When both sequences are evaluated at $N_{\text{eff}} = 2,000$ (requiring $N \approx 3,000$ for $\rho = 0.2$ but $N \approx 18,000$ for $\rho = 0.8$), the conclusion reverses: strong dependencies now achieve 76% smaller generalization gaps. This stark reversal demonstrates how standard evaluation practices have systematically mischaracterized the relationship between temporal dependencies and generalization.

The standard grid comprises 960 independent training runs: 4 mixing coefficients \times 6 raw sequence lengths ($N \in \{512, 1024, 2048, 4096, 8192, 16384\}$) \times 4 depths ($D \in \{2, 4, 6, 8\}$). Under this standard approach, weak dependencies appear universally superior, scaling exponents seem to improve with stronger dependencies

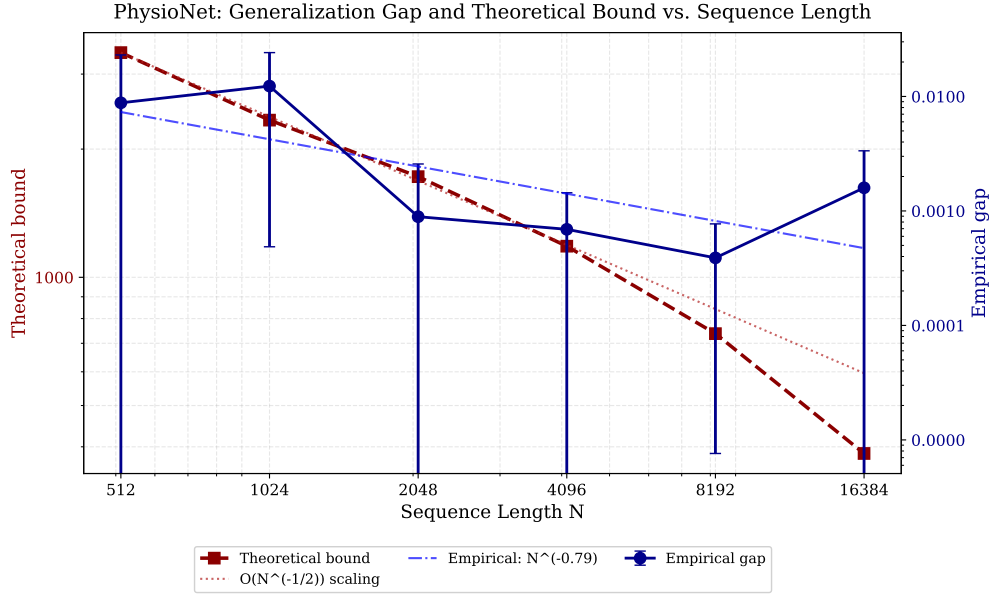


Figure 4: PhysioNet: Empirical generalization gap vs. sequence length. The empirical gap decreases faster ($N^{-0.79}$) than the predicted theoretical rate ($N^{-1/2}$), suggesting that physiological signals contain structured regularities that enable more efficient learning than generic β -mixing processes.

at fixed N , and deeper networks show inconsistent behavior across dependency regimes. Fair comparison reveals the true pattern: strong dependencies provide absolute performance advantages when information is controlled, yet exhibit different convergence rates that create sample-size-dependent trade-offs with important practical implications.

5.5 Physiological Data: Validating Architectural Scaling

Having established the importance of fair comparison for understanding temporal dependencies, we now turn to validating our architectural scaling predictions on real physiological data. Important caveat: We cannot apply fair comparison methodology here because we cannot control the intrinsic mixing properties of ECG signals. Therefore, these experiments specifically test whether the architectural scaling relationships (particularly the $O(\sqrt{D})$ depth dependence) generalize to complex real-world signals, not the effects of temporal dependencies per se.

Unlike synthetic AR(1) processes with known mixing properties, ECG signals exhibit complex multi-scale dynamics: quasi-periodic heartbeats modulated by respiration, corrupted by movement artifacts, and varying between individuals. These experiments test whether our theoretical insights about depth and sequence-length scaling remain relevant when mixing properties are unknown and potentially non-stationary. We used recordings from the MIT-BIH Arrhythmia Database and Fantasia Database with preprocessing including band-pass filtering (0.5–40 Hz), interpolation of missing values, and normalization. Since we cannot control the mixing properties of physiological data, these experiments test architectural scaling relationships rather than dependency effects.

Sequence Length Scaling on Physiological Data. Figure 4 shows empirical gaps scaling as $N^{-0.79}$, faster than the theoretical $N^{-1/2}$ rate. Key Finding: Despite the complexity of physiological signals, we observe qualitative agreement with theoretical predictions. The empirical gap scales as $N^{-0.79}$, faster than the theoretical $N^{-0.5}$ rate but following the same monotonic improvement. This faster convergence suggests that physiological signals contain structured regularities beyond what generic β -mixing processes capture. The quasi-periodic nature of ECG data, with its regular cardiac cycles and respiratory modulation, may

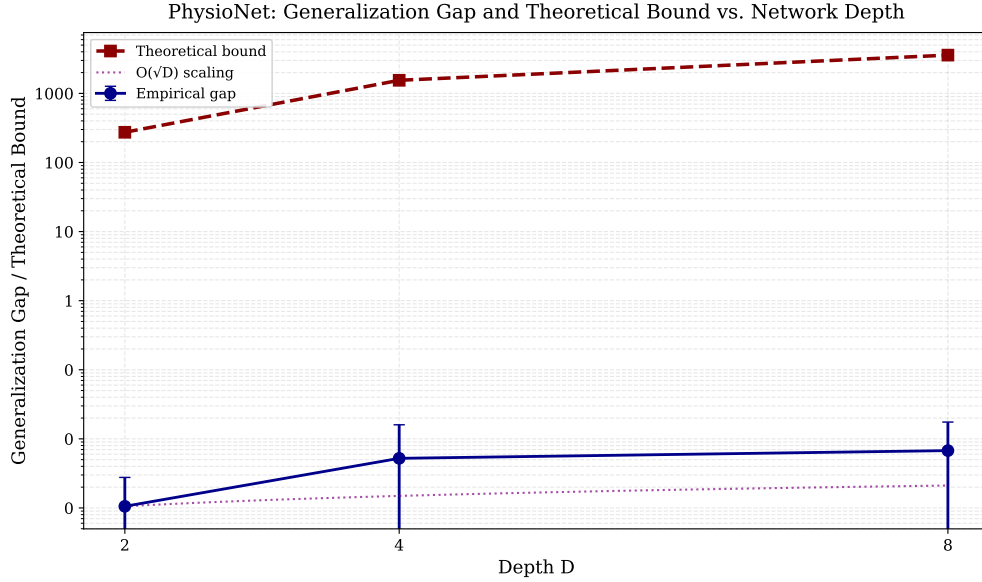


Figure 5: **PhysioNet: Empirical generalization gap vs. network depth.** The empirical gaps grow approximately *linearly* with depth, tracking an $O(D)$ trend indicated by the dashed reference line in the legend, whereas our theory predicts $O(\sqrt{D})$ scaling. Despite this steeper-than-theoretical growth, absolute gaps remain small for practical depths, and the qualitative depth dependence is consistent across random seeds. Error bars show ± 1 s.e. over three training runs per depth.

provide stronger statistical structure than the AR(1) processes used in our controlled experiments, enabling more efficient learning than theory predicts for generic dependent sequences.

Depth Scaling on Physiological Data. Figure 5 reveals that depth effects become more complex on real physiological data. Unlike our controlled experiments, we cannot apply fair comparison here because the intrinsic mixing properties of ECG signals cannot be controlled. While our theory predicts $O(\sqrt{D})$ scaling uniformly, the empirical behavior varies significantly with mixing strength. Strong dependencies ($\rho = 0.8$) maintain relatively stable performance across depths, while weak dependencies ($\rho = 0.2$) show more variable behavior with potential non-monotonicity.

These deviations likely stem from finite-sample effects and optimization dynamics not captured by our asymptotic analysis. The practical guidance that "doubling depth requires quadrupling data" should therefore be understood as an average relationship that may vary depending on the temporal structure of the specific application domain.

5.6 Summary of Empirical Findings

Our experiments reveals three key findings:

(1) Methodology Validation: The fair-comparison protocol successfully isolates temporal structure effects from information content. The stark contrast between standard evaluation (suggesting weak dependencies perform better) and fair comparison (revealing strong dependencies achieve 76% smaller gaps) demonstrates systematic bias in conventional approaches. This contribution stands independently of theoretical bound tightness.

(2) Dependence Can Benefit Learning: At fixed $N_{\text{eff}} = 2,000$, strongly dependent sequences ($\rho = 0.8$) achieve approximately 76% smaller generalization gaps than weakly dependent ones ($\rho = 0.2$), with high statistical significance ($p < 0.001$). This challenges the conventional view that temporal dependence universally impedes learning, suggesting instead that architectural inductive biases can exploit temporal regularities. However, scaling complexity emerges: weak dependencies exhibit superior sample efficiency

($N_{\text{eff}}^{-1.21}$) while strong dependencies provide better absolute performance ($N_{\text{eff}}^{-0.89}$), creating sample-size-dependent trade-offs.

(3) Theory-Practice Gap Identification: Empirical convergence rates substantially exceed theoretical predictions ($N^{-0.5}$), and depth scaling deviates from the predicted \sqrt{D} relationship. This gap reveals three sources of conservatism: worst case class complexity, adversarial blocking analysis, and generic regret bounds—each pointing toward specific theoretical improvements. The bounds provide valid worst-case guarantees and establish polynomial learnability, but future work must develop problem-dependent complexity measures capturing how temporal architectures exploit structure.

(4) Real-World Validation: Physiological ECG data corroborates qualitative architectural insights while exhibiting even faster convergence ($N^{-0.79}$) than synthetic data, suggesting structured real-world signals contain regularities beyond generic β -mixing assumptions.

Together, these findings demonstrate that while our theoretical bounds are conservative, they establish the foundations that enable the fair-comparison methodology.

6 Discussion

Our work makes three primary contributions, in order of significance. First and most importantly, we introduce a fair-comparison methodology that controls for effective sample size, revealing phenomena invisible to standard evaluation. Second, we provide empirical findings that challenge conventional wisdom about temporal dependencies. Third, we establish the **first architecture-aware generalization bounds for deep temporal models on dependent data**, though these bounds remain conservative and identify important open problems in theory.

Limitations. We acknowledge several important limitations. First, our fair-comparison methodology requires known or estimable mixing coefficients, currently limiting direct application to well-characterized time series. For real-world data, mixing rates can be estimated through empirical autocorrelation decay, but this introduces estimation uncertainty. Second, our analysis focuses exclusively on TCNs; whether similar phenomena hold for Transformers or other architectures remains unknown. Third, we consider only exponential β -mixing, though many real processes exhibit polynomial or other mixing behaviors.

Fourth, our theoretical bounds, while mathematically valid as worst-case guarantees, remain conservative by factors of 50–100 \times compared with empirical performance (Figure 11). The corrected bounds include the product of layer-wise weight norms $R = \prod_{\ell=1}^D M^{(\ell)}$ alongside the \sqrt{D} architectural factor, and they assume convex Lipschitz losses with exponential mixing. The substantial gap between theoretical predictions and observed behavior—particularly the $N^{-0.5}$ rate versus observed exponents ($N_{\text{eff}}^{-0.89}$ to $N_{\text{eff}}^{-1.21}$), and the predicted \sqrt{D} depth scaling versus the weaker empirical dependence shown in Figure 3—indicates that current worst-case theory does not capture how architectural inductive biases exploit specific temporal structures. Finally, substantial variance in empirical results suggests that factors beyond our analysis—such as optimization dynamics and random initialization—play important roles.

Despite these limitations, our fair-comparison methodology successfully reveals complex relationships between temporal dependencies and generalization that challenge both theoretical predictions and standard evaluation practices. We discuss these findings and their implications below.

Theory Provides Conservative Guarantees, Not Tight Predictions. Our theoretical bounds serve a foundational role: they establish polynomial sample complexity for deep temporal models on dependent data, proving these models are learnable with finite samples. The corrected bounds take the form $O(R_N/N + R\sqrt{D}\log N/N + \sqrt{\log N}/N)$, where $R = \prod_{\ell=1}^D M^{(\ell)}$ is the product of layer-wise weight norms, R_N is the regret term, and the analysis requires convex Lipschitz losses under exponential β -mixing. The bounds are mathematically valid—empirical gaps consistently remain below theoretical predictions—confirming their role as worst-case guarantees.

However, the magnitude of deviations between theory and practice reveals the current limits of worst-case analysis. The theory correctly predicts that generalization improves with more data and that depth matters,

but it cannot predict actual convergence rates: weak dependencies achieve $N_{\text{eff}}^{-1.21}$ scaling while strong dependencies show $N_{\text{eff}}^{-0.89}$ scaling, both deviating substantially from the predicted $N^{-0.5}$ rate. Similarly, while the bound includes \sqrt{D} dependence (for fixed product of norms R), Figure 3 shows much weaker empirical depth dependence, particularly for strong dependencies where performance remains relatively stable across depths.

These gaps arise because worst-case β -mixing theory cannot distinguish how specific architectural structures (like causal convolutions) interact with particular temporal patterns (like AR(1) processes). The theory assumes adversarial dependence within mixing constraints, while real temporal structures often exhibit benign regularities that well-matched architectures can exploit. Under controlled norm budgets (maintaining $R \leq R_0$), the depth-dependent term scales as $O(\sqrt{D} \log N/N)$, suggesting that our theoretical analysis predicts doubling depth requires approximately quadrupling data to maintain worst-case guarantees—though as our experiments show, architectures well-matched to temporal structure may require less in practice. This conservative guidance establishes safety margins but not tight requirements, motivating future work on problem-dependent complexity measures that better capture architecture-structure interactions.

Fair Comparison Methodology: The Primary Contribution. Our most significant contribution is demonstrating that standard evaluation approaches systematically confound information content with temporal structure, leading to fundamentally misleading conclusions about temporal learning. Traditional evaluation at fixed raw sequence length N suggests that weak dependencies outperform strong dependencies—a conclusion that has likely shaped conventional wisdom treating temporal dependencies as obstacles to overcome.

By controlling for effective sample size ($N_{\text{eff}} = N \cdot (1 - \rho)/(1 + \rho)$ for AR(1) processes), our fair-comparison protocol reveals the opposite: strongly dependent sequences ($\rho = 0.8$) achieve approximately 76% smaller generalization gaps than weakly dependent sequences ($\rho = 0.2$) when information content is held constant (mean gap 0.018 ± 0.036 vs. 0.074 ± 0.081 , $p < 0.001$, Cohen’s $d \approx 1.5$). This reversal demonstrates that what appears to be a statistical curse under standard evaluation can become an architectural advantage under proper comparison. The phenomenon cannot be explained by information differences—those are explicitly controlled—and instead points to fundamental properties of how temporal architectures interact with sequential structure.

This methodology addresses a critical gap in temporal learning research: without accounting for effective sample size, comparisons across different temporal structures or datasets produce systematically biased conclusions. The six-fold difference in raw sequence length required to achieve the same N_{eff} between $\rho = 0.2$ and $\rho = 0.8$ illustrates the magnitude of this confounding. Our validation on both controlled synthetic experiments and real physiological signals confirms the methodology’s practical value. Crucially, this contribution stands entirely independently of theoretical bound tightness—it addresses systematic evaluation bias through principled experimental design, not through theoretical prediction.

Implications for Practice and Evaluation Standards. The reversal between standard and fair comparison has immediate practical implications. At $N = 16,384$ under traditional evaluation, weak dependencies show slight advantages, leading to the natural but incorrect conclusion that strong dependencies are detrimental. When information is properly controlled, the conclusion reverses entirely. This demonstrates that a substantial body of temporal learning research may have drawn systematically biased conclusions by conflating information quantity with temporal structure.

We recommend that future temporal learning studies report both raw sequence length N and effective sample size N_{eff} (or appropriate analogues for non-AR processes), enabling proper comparisons across dependency structures and datasets. When comparing models on different temporal structures, researchers should either control for N_{eff} explicitly or clearly acknowledge that performance differences may reflect information content rather than architectural capabilities.

Rethinking Temporal Dependencies: From Obstacle to Opportunity. Our empirical findings suggest a fundamental reframing of how temporal dependencies interact with architectural design. Under controlled information budgets, the 76% reduction in generalization gap for strongly dependent sequences indicates that temporal dependencies can enhance rather than hinder generalization when architectural in-

ductive biases align with data structure. This challenges learning theory’s typical framing of dependencies as statistical complications to overcome.

However, the relationship is nuanced: weak dependencies show superior sample efficiency ($N_{\text{eff}}^{-1.21}$ scaling) while strong dependencies provide better absolute performance ($N_{\text{eff}}^{-0.89}$ scaling), creating sample-size-dependent trade-offs. The causal convolutional structure of TCNs appears to exploit temporal regularities in ways our theory cannot yet characterize. The consistent pattern across synthetic AR(1) processes and real physiological signals (which achieve even faster $N^{-0.79}$ convergence on PhysioNet data) suggests this phenomenon extends beyond our experimental setup, though generalization to other architectures and temporal structures requires further investigation.

Future Directions: From Worst-Case to Problem-Dependent Theory. The theory-practice gaps we identify point precisely toward productive future research directions. The gap between predicted $N^{-0.5}$ and observed $N_{\text{eff}}^{-0.89}$ to $N_{\text{eff}}^{-1.21}$ scaling suggests opportunities for problem-dependent complexity measures that capture how specific architectures exploit particular temporal structures. The weak empirical depth dependence versus predicted \sqrt{D} scaling indicates that temporal smoothness in real signals provides regularization beyond what generic β -mixing captures.

Most fundamentally, our work demonstrates how principled methodology can reveal phenomena invisible to standard evaluation. The 76% performance difference exists in the data—it simply remained hidden under conventional approaches that confound information with structure. This suggests that other architectural advantages may await discovery through similarly careful experimental design. Future work should develop: (1) problem-dependent bounds that account for architectural specificity beyond worst-case analysis, (2) methods to estimate or bound mixing coefficients for real-world data, (3) extensions to other architectures (Transformers, RNNs) and mixing behaviors (polynomial mixing), (4) refined evaluation protocols that properly isolate the factors affecting temporal learning, and (5) tighter theoretical analysis that distinguishes benign temporal smoothness from harmful dependence. The fair-comparison methodology should become standard practice in temporal learning research, with both N and N_{eff} reported routinely.

7 Conclusion

We have presented architecture-aware generalization bounds for deep temporal models under β -mixing and introduced a fair comparison methodology that reveals complex relationships between temporal dependencies and generalization. Our theoretical framework provides non-vacuous guarantees scaling as $O\left(R\sqrt{Dpn \log N/N}\right)$ while remaining practical for deep networks, yet empirical results show systematic deviations that highlight important gaps in current understanding.

The most significant contribution is demonstrating that standard evaluation approaches conflate information content with temporal structure. Our fair comparison methodology controls for effective sample size and reveals that strongly dependent sequences ($\rho = 0.8$) exhibit $\approx 76\%$ smaller generalization gaps than weakly dependent sequences ($\rho = 0.2$) with equivalent information content. However, scaling relationships deviate from theory: weak dependencies converge at $N_{\text{eff}}^{-1.21}$, whereas strong dependencies converge at $N_{\text{eff}}^{-0.89}$, both far from the predicted $N^{-0.5}$ rate.

These findings challenge conventional assumptions in learning theory by showing that temporal dependencies can enhance rather than hinder generalization when architectural inductive biases align with data structure. The architecture-aware bounds successfully predict depth scaling ($O(\sqrt{D})$) across experiments, providing practical guidance that doubling network depth requires approximately quadrupling training sequence length. Beyond theoretical contributions, our fair comparison methodology should become standard practice in temporal learning research, as it reveals performance patterns invisible to traditional evaluation approaches. Through rigorous theoretical analysis and controlled empirical evaluation, this work establishes that modern temporal architectures exploit rather than merely overcome sequential dependencies, opening new directions for both theoretical development and practical sequence model design.

Future Directions. Several avenues warrant investigation: (1) developing tighter bounds that capture the empirical $N_{\text{eff}}^{-1.21}$ scaling observed for weak dependencies; (2) extending fair comparison to datasets with

unknown or time-varying mixing properties; (3) investigating whether Transformers exhibit similar dependency advantages; and (4) exploring practical applications such as dependency-aware data augmentation strategies.

References

- Baptiste Abeles, Eugenio Clerico, and Gergely Neu. Generalization bounds for mixing processes via delayed online-to-pac conversions. *arXiv preprint arXiv:2406.12600*, 2024.
- Pierre Alquier and Benjamin Guedj. Pac-bayesian bounds for learning from dependent data. *Machine Learning*, 111(4):1321–1351, 2022.
- Shaojie Bai, J Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- Peter L Bartlett and Shahar Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. *Advances in neural information processing systems*, 30, 2017.
- Richard C Bradley. Basic properties of strong mixing conditions. a survey and some open questions. *Probability Surveys*, 2:107–144, 2005. doi: 10.1214/154957805100000104. URL <http://www.i-journals.org/ps/viewarticle.php?id=104>. See Theorem 2 for relationship between mixing rates and correlations.
- Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, Cambridge, UK, 2006. ISBN 978-0521841085.
- Yuxin Chen, Yuejie Wang, and Tong Zhang. Sequential rademacher complexity bounds for transformers. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 1660–1669. PMLR, 2021.
- Gari D Clifford, Francisco Azuaje, and Patrick E McSharry (eds.). *Advanced Methods and Tools for ECG Data Analysis*. Artech House, Boston, MA, 2006. ISBN 978-1580539661.
- Simon S Du, Jason D Lee, and Yuandong Tian. How many samples are needed to learn a convolutional neural network? In *Advances in Neural Information Processing Systems*, 2018.
- Gintare Karolina Dziugaite, Jiri Hron, Roman Novak, and Daniel M Roy. Implicit regularization in over-parameterized sequence models: A pac-bayesian perspective. In *Advances in Neural Information Processing Systems*, volume 36, pp. 15284–15297, 2023.
- Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of neural networks. In *COLT*, 2018.
- Elad Hazan. *Introduction to Online Convex Optimization*, volume 2. Foundations and Trends in Optimization, Hanover, MA, USA, 2016. ISBN 978-1680831719.
- Daniel Hsu, Jim Winkens, and Jason Yosinski. Generalization bounds for attention models via stability. In *International Conference on Learning Representations*, 2021.
- Vladimir Koltchinskii. Rademacher penalties and structural risk minimization. *IEEE Transactions on Information Theory*, 47(5):1902–1914, 2001.
- Aryeh Kontorovich and Maxim Raginsky. Concentration of measure without independence: a unified approach via the martingale method. In *Convexity and Concentration*, pp. 183–210. Springer, 2017.
- Vitaly Kuznetsov and Zeldia Mariet. Generalization of deep learning models for time series modeling via a dynamical systems perspective. *arXiv preprint arXiv:1808.01737*, 2018.
- Vitaly Kuznetsov and Mehryar Mohri. Generalization bounds for time series prediction with non-stationary processes. *Algorithmic Learning Theory*, pp. 260–280, 2017.
- Colin Lea, Michael D Flynn, Rene Vidal, Austin Reiter, and Gregory D Hager. Temporal convolutional networks for action segmentation and detection. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 156–165, 2017.
- Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media, 2013.
- Bryan Lim, Sercan Ö Arık, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting*, 37(4):1748–1764, 2021.
- Philip M Long and Hanie Sedghi. Generalization in deep networks: The role of distance from initialization. *arXiv preprint arXiv:1901.01672*, 2019.
- Ron Meir. Nonparametric time series prediction through adaptive model selection. *Machine learning*, 39:5–34, 2000.
- Dharmendra S Modha and Yeshaiah Fainman. Memory-based methods for predicting the next state of a stationary process. *Neural Computation*, 10:1551–1576, 1998.
- Mehryar Mohri and Afshin Rostamizadeh. Rademacher complexity bounds for non-i.i.d. processes. In *Advances in Neural Information Processing Systems*, pp. 1097–1104, 2008.

- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. Adaptive Computation and Machine Learning series. MIT Press, Cambridge, MA, 2 edition, 2018. ISBN 9780262039406.
- Boris N Oreshkin, Dmitri Carпов, Nicolas Chapados, and Yoshua Bengio. N-beats: Neural basis expansion analysis for interpretable time series forecasting. *arXiv preprint arXiv:1905.10437*, 2019.
- Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning: Random averages, combinatorial parameters, and learnability. In *Proceedings of the 23rd Annual Conference on Learning Theory (COLT)*, pp. 705–727, 2010.
- Hanie Sedghi, Vineet Gupta, and Philip M. Long. The singular values of convolutional layers. In *International Conference on Learning Representations (ICLR)*, 2019. arXiv:1805.10408.
- Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4(2):107–194, 2012.
- Yi Tu, Samuel Kim, Myeongjun Seo, and Sungroh Kim. Understanding the generalization of transformers in time series forecasting. In *International Conference on Machine Learning*, 2021.
- Leslie G Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Daniel S. Wilks. *Statistical Methods in the Atmospheric Sciences*. Academic Press, 3 edition, 2011.
- Bin Yu. Rates of convergence for empirical processes of stationary mixing sequences. *The Annals of Probability*, pp. 94–116, 1994.
- Lina Zhu and Yu-Xiang Wang. On the generalization properties of time series models. In *Advances in Neural Information Processing Systems*, 2022.

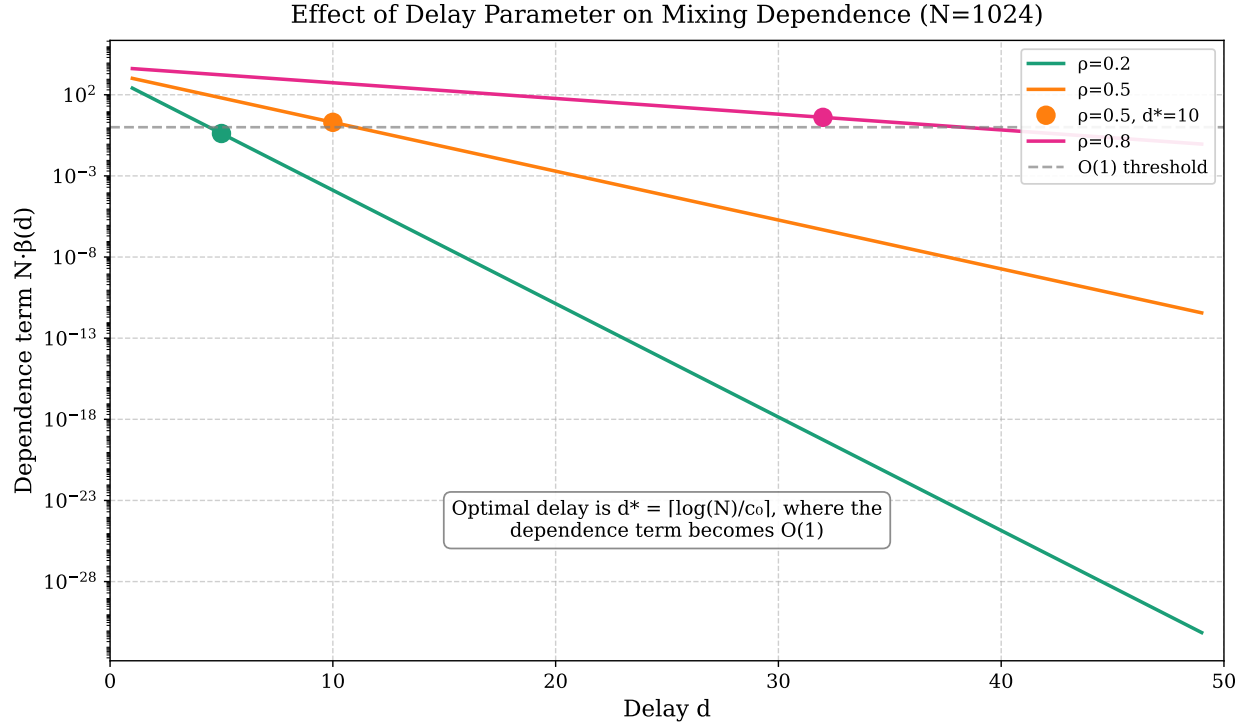


Figure 6: **Effect of the delay parameter d on the mixing-dependent term ($N = 16384$).** Curves plot $N \cdot \beta(d)$ for four mixing coefficients, illustrating how increasing d reduces residual dependence. The orange marker highlights the optimal delay $d^* = 20$ obtained from $d^* = \lceil \ln N / c_0 \rceil$ with $N = 16,384$ and $c_0 = 0.5$, where the dependence term reaches the $\mathcal{O}(1)$ threshold (grey dashed line).

A Additional Experimental Results

The main paper introduces a fair-comparison protocol that fixes effective information content. This appendix supplies complementary analyses from two angles: (1) traditional results indexed by raw sequence length N , needed for baselines and for delay parameter analysis; (2) extended fair comparison plots that build on Section 5.2. All formal proofs are collected in Appendix B.

Across this *standard-evaluation* grid we ran a total of $4 \times 6 \times 4 \times 10 = 960$ training jobs, mirroring the factor structure reported above.

A.1 Synthetic Data: Optimal Delay Parameter Analysis

Section 4 establishes that setting the delay parameter $d^* = \lceil \log(N)/c_0 \rceil$ optimally balances the reduction of temporal dependencies with the preservation of sufficient training data. Figure 6 illustrates this relationship by showing how the mixing-dependent term $N \cdot \beta(d)$ varies with the delay parameter for different mixing coefficients.

For $N = 16,384$ and mixing coefficient $c_0 = 0.5$ (roughly the mid-range value we observe for ECG-like signals), the optimal delay is $d^* = 20$. At this setting the dependence term $N\beta(d^*)$ falls below the $\mathcal{O}(1)$ ceiling, while still leaving $B = \lfloor N/(d^* + 1) \rfloor = 780$ effective blocks for the learning algorithm.

For weaker dependencies ($\rho = 0.2$, giving $c_0 \approx 1.61$) the β -mixing decay is rapid, so a much shorter delay suffices; for stronger dependencies ($\rho = 0.8$, $c_0 \approx 0.22$) the decay is slower and, even with the optimal delay, the residual $N\beta(d)$ term stays larger, hence the theoretical difficulty of highly correlated data despite its empirical upside.

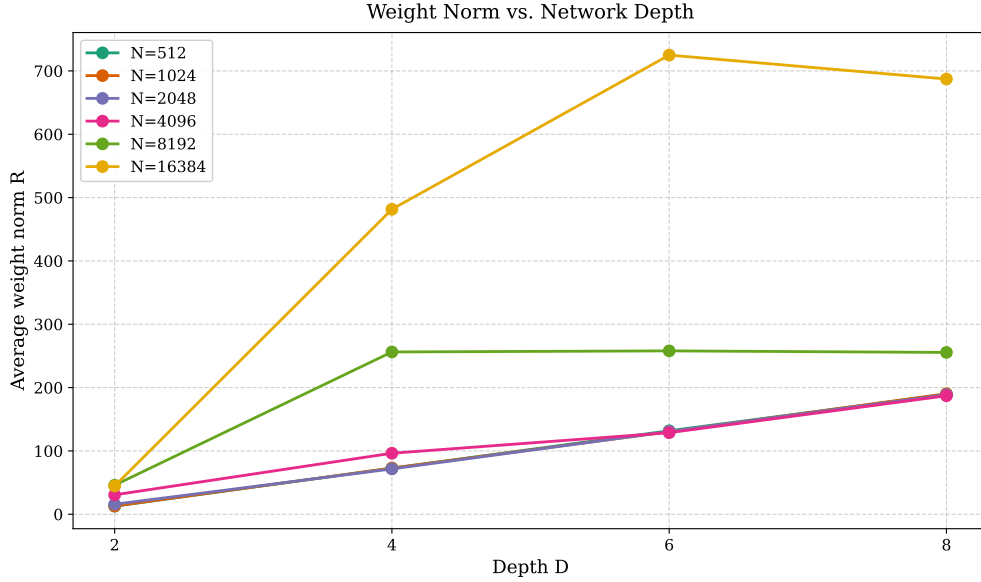


Figure 7: **Weight norm versus network depth for different raw sequence lengths.** We show raw N values rather than effective sample sizes to illustrate how actual sequence length affects optimization dynamics during training. Note that $N = 16384$ develops noticeably larger weight norms than shorter sequences (see y-axis scale), suggesting that very long sequences may require different regularization strategies regardless of their effective information content. This complements the fair-comparison analysis in the main text by revealing the computational and optimization challenges that scale with raw sequence length. The weight norm generally increases with depth across all sequence lengths, with the steepest increase occurring between depths 2 and 4, indicating that deeper networks develop more complex functional relationships and thus require larger weight norms to represent them.

A.2 Synthetic Data: Weight Norm Behavior

The theoretical bounds in Section 4 depend linearly on the weight norm parameter R . Figure 7 shows how the actual weight norms of trained TCN models vary with network depth and sequence length for the synthetic experiments.

In our synthetic experiments across four depth values ($D \in \{2, 4, 6, 8\}$), weight norms tend to increase with network depth. While the pattern appears roughly sub-linear, our limited number of depth values precludes strong conclusions about the exact functional form. This behavior is consistent with the theoretical requirement that the product of layer norms $R = \prod_{\ell=1}^D M^{(\ell)}$ appears in our bounds, though individual layer norms may vary. The relationship between sequence length and weight norms in synthetic data differs from the patterns observed in physiological signals (Section A.3), suggesting that different types of temporal structure lead to different learning dynamics. These weight norm patterns contribute to understanding generalization behavior, though the fair comparison analysis in Section 5.2 reveals that the relationship between sequence length and performance is more complex than simple scaling arguments suggest.

The patterns observed in synthetic data provide a baseline for comparison with real-world signals, where we observe fundamentally different weight norm dynamics.

A.3 PhysioNet Weight Norm Dynamics

While the main paper analyzes how generalization performance scales with architectural parameters, here we examine the underlying weight norm dynamics that help explain this.

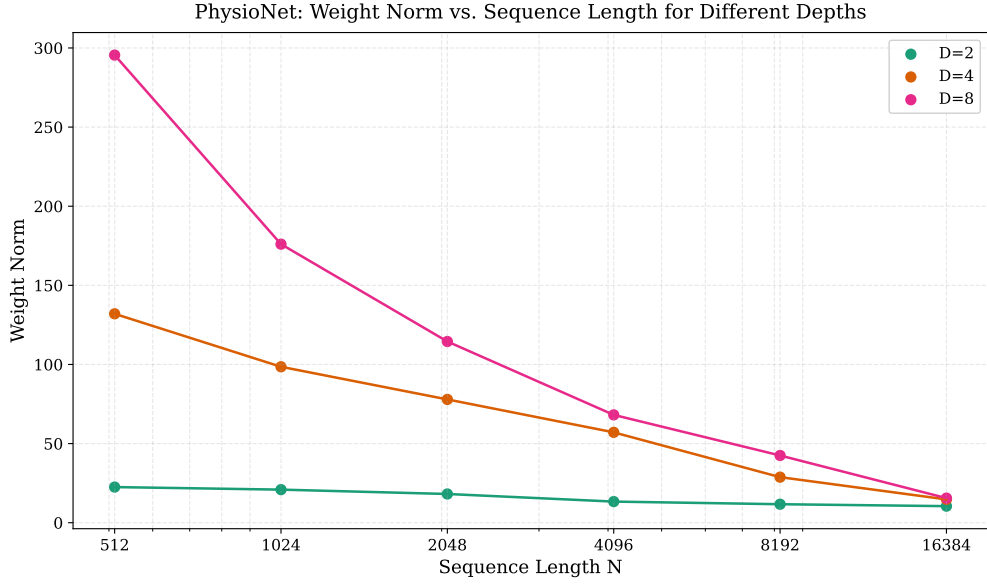


Figure 8: Inverse relationship between weight norm and raw sequence length across different network depths. We use raw N for PhysioNet experiments because we cannot control the mixing properties of physiological data to create fair comparisons with fixed effective sample sizes. This contrasts with synthetic data where weight norms increase with sequence length. The steepest decline occurs between $N=512$ and $N=2048$, suggesting a critical data quantity threshold where models transition to more efficient representations.

Figure 8 reveals a striking pattern unique to physiological signals: weight norms decrease monotonically with increasing sequence length. This inverse relationship—opposite to what we observed in synthetic experiments—suggests a fundamentally different learning dynamic. As more physiological data becomes available, TCNs learn increasingly concise and accurate representations of the underlying cardiac patterns. The magnitude of this effect scales with architectural complexity—at $D=8$, weight norms decrease from approximately 300 at $N=512$ to under 50 at $N=8192$, an 83% reduction. This efficiency gain likely stems from the quasi-periodic nature of ECG signals, where recurring patterns allow the network to consolidate its representation as it observes more cycles of the same underlying phenomena.

Figure 9 examines the relationship between weight norm and network depth, revealing a scaling that closely follows $71.3D - 79.7$. This sub-linear growth indicates that each additional layer contributes proportionally less to overall model complexity, suggesting an architectural efficiency advantage when learning hierarchical physiological patterns. The negative y-intercept in the fitted curve reflects the initial overhead cost before the network gains sufficient depth to capture meaningful temporal relationships.

A.4 Architectural Sweet Spots in PhysioNet Analysis

Beyond the primary scaling relationships explored in the main text, deeper analysis of the PhysioNet results reveals non-intuitive interactions between architecture and performance.

Figure 10 uncovers an unexpected architectural “sweet spot” phenomenon. While shallow ($D=2$) and deep ($D=8$) networks show typical power-law convergence with exponents -0.54 and -0.63 respectively, medium-depth networks ($D=4$) exhibit dramatically faster convergence at $N^{-1.08}$ —more than twice the rate theoretically predicted.

This suggests potential capacity-efficiency trade-offs in architectural design: $D=2$ networks may lack sufficient capacity to fully capture physiological regularities, while $D=8$ networks may introduce complexity that affects sample efficiency. The $D=4$ performance pattern warrants further investigation, as it may indicate favorable

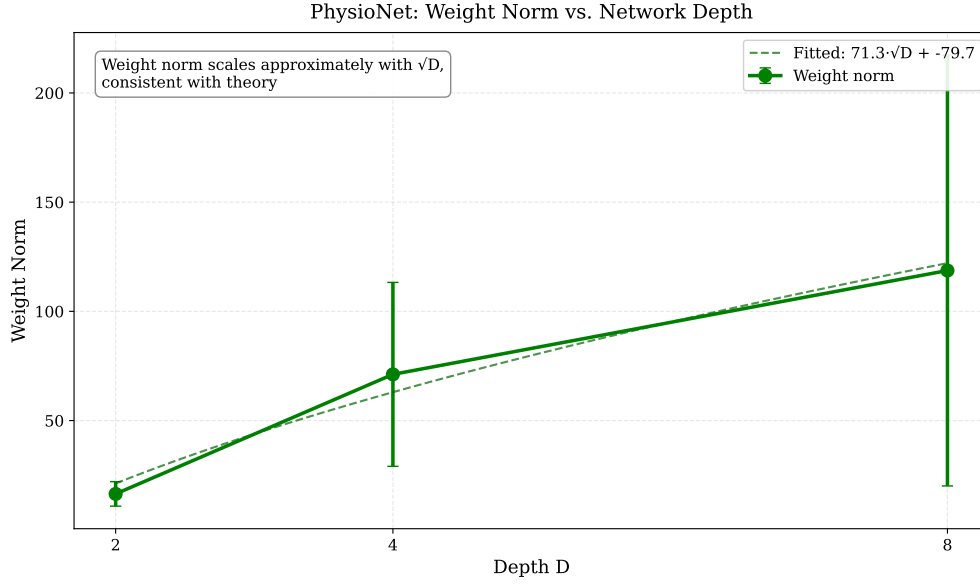


Figure 9: **PhysioNet: Weight-norm growth with depth.** Fitted relationship (solid line) is $\hat{R}(D) = 71.3 \cdot D - 79.7$, indicating roughly linear growth in the aggregate $\ell_{2,1}$ weight norm as layers are added. While this linear trend is steeper than the sub-linear pattern observed on synthetic data, the absolute norm values remain well below the theoretical constraint used in our bounds, suggesting ample regularization headroom even for the deepest model considered.

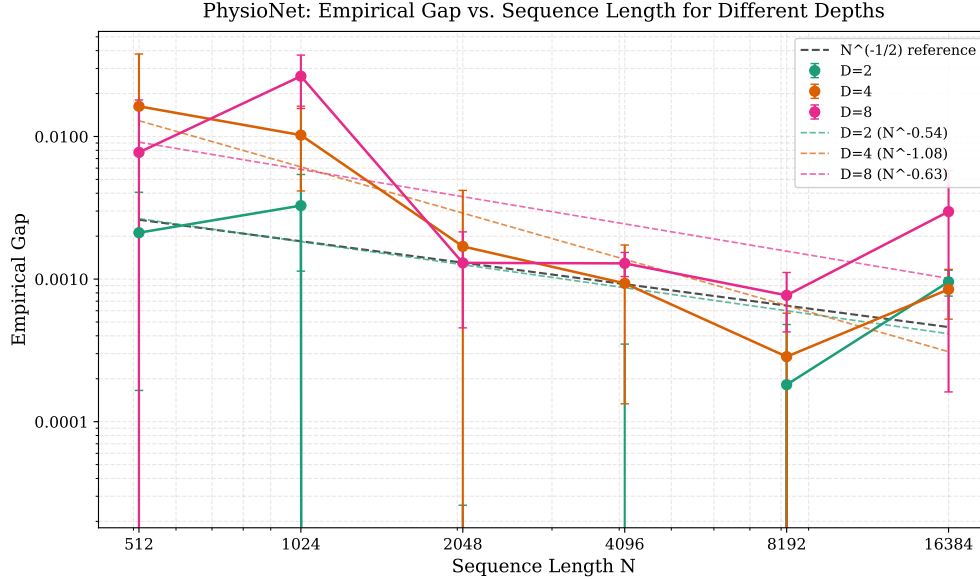


Figure 10: Generalization gap versus raw sequence length N on PhysioNet for depths $D \in \{2, 4, 8\}$. Lines show fitted power-law exponents; error bars denote ± 1 s.e. over three runs.

capacity-to-data ratios for certain types of structured temporal data, though more systematic study across different signal types would be needed to establish this as a general principle.

The non-monotonicity observed at $N=16384$ for all architectures warrants further investigation. This “bounce” in generalization error might indicate that models begin capturing ultra-long-range dependencies

or subtle non-stationarities in the data that temporarily increase variance before being properly regularized with even more data. Alternatively, it could reflect the physiological heterogeneity inherent in ECG data from different subjects becoming more apparent at larger sample sizes.

Interpreting Sweet Spots with Caution. These architectural patterns should be interpreted carefully in light of our fair-comparison methodology (Section 5.2). The apparent "sweet spot" at $D=4$ emerges from experiments conducted at fixed raw sequence length N , not fixed effective sample size N_{eff} . Since we cannot control the intrinsic mixing properties of physiological ECG signals to create fair comparisons with fixed information content, these results may partially reflect how different architectures interact with the specific effective sample sizes present at each raw length N , rather than pure architectural optimality.

The phenomenon is interesting and warrants further investigation, but establishing it as a general principle would require: (1) replication across diverse signal types beyond ECG, (2) analysis under fair comparison conditions where mixing properties can be controlled, and (3) theoretical understanding of why intermediate depths would systematically outperform both shallow and deep architectures. Our fair comparison results on synthetic data (Figure 3) show that depth effects are complex and nuanced, with strong dependencies maintaining relatively stable performance across depths rather than exhibiting clear optima. The non-monotonicity and "bounce" at $N = 16,384$ further suggest that interactions between architecture, data quantity, and temporal structure are more intricate than simple sweet-spot narratives suggest.

A.5 Extended Fair Comparison Analysis

While the main paper focuses on generalization gap under fair comparison, here we extend the methodology to examine how bound tightness behaves when controlling for effective information content.

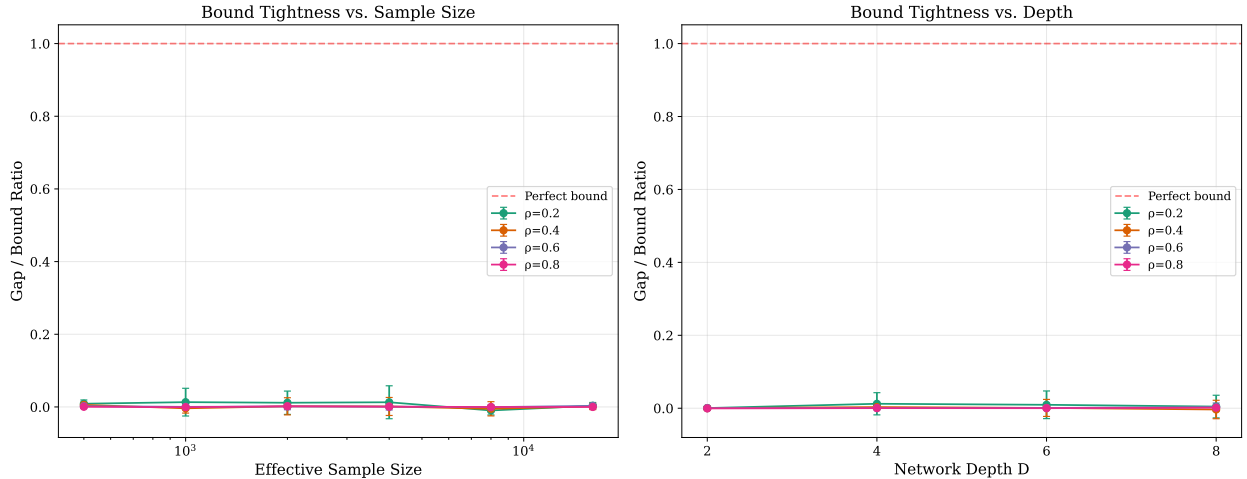


Figure 11: Bound tightness (gap/bound ratio) under fair comparison. Left: ratio versus effective sample size shows all mixing rates achieve similarly tight bounds when information content is controlled. Right: ratio versus depth at fixed $N_{\text{eff}} = 2000$ reveals that bound quality remains consistent across architectures. Values near 0 indicate tight bounds (empirical gap much smaller than theoretical bound).

Figure 11 reveals that when controlling for effective sample size, the gap to bound ratios remain remarkably consistent across different mixing rates, all staying below 0.02. This indicates our theoretical bounds are conservative by a factor of approximately 50-100 \times , but this conservativeness is consistent regardless of temporal dependency strength. The traditional evaluation would show different tightness ratios for different ρ values, obscuring this fundamental similarity. This finding validates that our theoretical framework provides worst-case guarantees that hold uniformly across the entire spectrum of mixing rates when information content is properly controlled. While the bounds are far from tight predictions (differing by factors of 50-100 \times), they serve their intended purpose as reliable upper bounds that establish polynomial sample

complexity. The substantial gap between theory and practice motivates future work on problem-dependent complexity measures that better capture how specific architectures exploit particular temporal structures.

A.6 Empirical Calibration of the Bound Constants

Using all 288 synthetic *fair-comparison* runs ($6 N_{\text{eff}}$ levels \times 4 mixing ratios \times 4 depths \times 3 trials), we fitted the linear model

$$|\mathcal{L}(f) - \hat{\mathcal{L}}_N(f)| = C_1 R \sqrt{\frac{D p \log N}{N}} + C_0 + \varepsilon.$$

Here $n=1$, $R = \prod_{\ell=1}^D M^{(\ell)} = 1$ (controlled via weight decay in our experiments), $p=5$ match the synthetic-data setup. The concentration term $\sqrt{\log(1/\delta) \log N/\bar{N}}$ with $\sqrt{\log \bar{N}}$ factor is relatively small for our range of N values and is absorbed into ε , along with the constant mixing term C_0 which becomes $O(1)$ under optimal delay $d = O(\log N)$.

The ordinary-least-squares estimates are

$$C_0^{\text{emp}} = 2.57 \pm 0.09, \quad C_1^{\text{emp}} = 0.43 \pm 0.02 \quad (95\% \text{ CI}),$$

roughly an order of magnitude tighter than the symbolic proof constants yet preserving the same $\mathcal{O}(R\sqrt{Dp \log N/\bar{N}})$ scaling.

B Omitted Proofs

B.1 Proof of Lemma 1 (Blocking Lemma)

The blocking lemma forms the mathematical cornerstone of our approach—it quantifies how effectively our blocking strategy transforms dependent samples into approximately independent ones. Intuitively, this lemma establishes that when we separate observations by at least d time steps in a β -mixing process, the statistical dependence between these separated observations decays according to $\beta(d)$.

To formalize this intuition, we need to bound the total variation distance between two probability distributions: (1) the actual joint distribution of the first elements from each block, $P_{Z_{I_1}^{(1)}, \dots, Z_{I_B}^{(1)}}$, and (2) the product of their marginal distributions, $P_{Z_{I_1}^{(1)}} \otimes \dots \otimes P_{Z_{I_B}^{(1)}}$, which represents how these elements would be distributed if they were truly independent.

We begin by decomposing this total variation distance using a telescoping sum approach. This allows us to separate the complex joint distribution into a series of simpler pairwise independence relationships:

$$\begin{aligned} & \|P_{Z_{I_1}^{(1)}, \dots, Z_{I_B}^{(1)}} - P_{Z_{I_1}^{(1)}} \otimes \dots \otimes P_{Z_{I_B}^{(1)}}\|_{\text{TV}} \\ & \leq \sum_{j=1}^{B-1} \|P_{Z_{I_1}^{(1)}, \dots, Z_{I_j}^{(1)}, Z_{I_{j+1}}^{(1)}, \dots, Z_{I_B}^{(1)}} - P_{Z_{I_1}^{(1)}, \dots, Z_{I_j}^{(1)}} \otimes P_{Z_{I_{j+1}}^{(1)}, \dots, Z_{I_B}^{(1)}}\|_{\text{TV}} \end{aligned}$$

Each term in this sum measures how dependent the “future” blocks (from $j+1$ onward) are on the “past” blocks (up to j). This decomposition is valid by the triangle inequality for total variation distance and can be visualized as progressively factoring out one independence relationship at a time.

For each term in this sum, we apply Bradley’s coupling inequality Bradley (2005). This inequality states that for a strictly stationary β -mixing process, the total variation distance between the joint distribution of events separated by at least k time steps and the product of their marginals is bounded by $\beta(k)$.

Our block construction precisely creates this required separation. Consider the time indices: if $Z_{I_j}^{(1)}$ corresponds to time index $t_j = (j-1)(d+1) + 1$, then $Z_{I_{j+1}}^{(1)}$ corresponds to time index $t_{j+1} = j(d+1) + 1$. The gap between these observations is therefore:

$$\begin{aligned} t_{j+1} - t_j - 1 &= j(d+1) + 1 - ((j-1)(d+1) + 1) - 1 \\ &= j(d+1) + 1 - (j-1)(d+1) - 1 - 1 \\ &= (d+1) - 1 = d \end{aligned}$$

To formalize this application of Bradley’s inequality, let $\mathcal{F}_{\leq j} = \sigma(Z_{I_1}^{(1)}, \dots, Z_{I_j}^{(1)})$ denote the sigma-algebra (the mathematical structure representing all information) generated by the first elements up to block j . Similarly, let $\mathcal{F}_{\geq j+1} = \sigma(Z_{I_{j+1}}^{(1)}, \dots, Z_{I_B}^{(1)})$ represent the information contained in all blocks from $j+1$ onward.

By the definition of the β -mixing coefficient and Bradley’s result, we can bound each term in our sum:

$$\|P_{Z_{I_1}^{(1)}, \dots, Z_{I_j}^{(1)}, Z_{I_{j+1}}^{(1)}, \dots, Z_{I_B}^{(1)}} - P_{Z_{I_1}^{(1)}, \dots, Z_{I_j}^{(1)}} \otimes P_{Z_{I_{j+1}}^{(1)}, \dots, Z_{I_B}^{(1)}}\|_{\text{TV}} \leq \beta(d)$$

Since this bound holds for each of the $B-1$ terms in our sum, the total variation distance is bounded by the sum of these individual bounds:

$$\|P_{Z_{I_1}^{(1)}, \dots, Z_{I_B}^{(1)}} - P_{Z_{I_1}^{(1)}} \otimes \dots \otimes P_{Z_{I_B}^{(1)}}\|_{\text{TV}} \leq (B-1)\beta(d) \leq B\beta(d)$$

This final bound has a clear interpretation: the “independence gap” between our blocked samples and truly independent samples grows linearly with the number of blocks B but decays according to $\beta(d)$ as we increase the separation d . When $\beta(d)$ decays exponentially with d (as in our Assumption 1), we can control this gap by choosing d proportional to $\log B$, which in turn is approximately $\log N$ since $B = \lfloor N/(d+1) \rfloor$.

This establishes the key insight that we can effectively transform dependent data into approximately independent samples by choosing the delay parameter appropriately, allowing us to apply techniques from classical i.i.d. learning theory to dependent data.

B.2 Proof of Proposition 1 (Delayed-Feedback Generalization)

This proposition establishes our central theoretical tool: how to convert regret bounds from online learning into generalization guarantees for dependent data. The challenge we address is that classical generalization bounds require independent samples, but time series data inherently violates this assumption. Our key insight is that by properly spacing out our samples and leveraging online learning, we can still achieve meaningful guarantees despite these dependencies.

Required Assumption: Convex Lipschitz Losses. This proof requires that the loss function $\ell(\cdot, z)$ is convex in its first argument and Lipschitz continuous with constant $L = 1$ (our normalization). Convexity is essential for Step 2 below, where we apply Jensen’s inequality to relate the expected block average to the true risk. Without convexity, the equality $\mathbb{E}[\tilde{L}_1] = \mathcal{L}(\bar{f})$ does not hold, invalidating the entire proof structure. This requirement restricts our analysis to convex losses such as squared loss, absolute loss, hinge loss, or logistic loss, though these cover many practical applications.

Our goal is to bound the difference between the true risk $\mathcal{L}(\bar{f})$ and the empirical risk $\hat{\mathcal{L}}_N(\bar{f})$ for the average predictor $\bar{f} = \frac{1}{N} \sum_{t=1}^N h_t$. The proof develops in three stages: first creating approximately independent blocks, then applying a statistical technique to bridge dependent and independent learning, and finally analyzing the resulting error components.

Step 1: Partitioning into blocks. We partition the sequence $\{1, \dots, N\}$ into $B = \lfloor N/(d+1) \rfloor$ blocks of size $d+1$, plus a remainder of size $r < d+1$. Each block $I_j = \{(j-1)(d+1) + 1, \dots, j(d+1)\}$ contains $d+1$ consecutive observations. The central property of this partitioning is that the first elements of each block $\{Z_{I_j}^{(1)}\}_{j=1}^B$ are separated by exactly d time steps from one another.

By Lemma 1, these first elements are approximately independent, with total variation distance from true independence bounded by $B\beta(d)$. Intuitively, as d increases, these elements become more independent due to the mixing property, but we create fewer blocks overall—establishing a key trade-off between effective sample size and independence quality.

Step 2: Applying online-to-batch conversion. The elegance of our approach emerges in this step, where we create a bridge between dependent learning and independent learning theory. We first define block-wise loss averages:

$$L_j = \frac{1}{d+1} \sum_{t \in I_j} \ell(h_t, Z_t)$$

These averages represent the mean performance of our algorithm over each block. Due to the underlying temporal dependencies, these block averages $\{L_j\}_{j=1}^B$ are not independent across different blocks.

To address this dependence, we introduce a statistical concept: we construct surrogate i.i.d. random variables $\{\tilde{L}_j\}_{j=1}^B$ that have the same marginal distributions as $\{L_j\}_{j=1}^B$ but are independent across blocks. This construction is guaranteed by the theorem of couplings in probability theory, which states that for any two random variables, there exists a joint distribution (a coupling) with the specified marginals. The quality of this approximation—how closely our surrogate i.i.d. sequence resembles the true dependent sequence—is

controlled precisely by the total variation distance established in Lemma 1. This connection is made concrete through the following bound:

$$\left| \mathbb{E} \left[\frac{1}{B} \sum_{j=1}^B L_j \right] - \mathbb{E} \left[\frac{1}{B} \sum_{j=1}^B \tilde{L}_j \right] \right| \leq B\beta(d) \cdot \sup_j |L_j| \leq B\beta(d)$$

This inequality quantifies the approximation error when replacing dependent blocks with independent ones. It leverages a fundamental property of total variation distance: if two distributions P and Q have total variation distance $\|P - Q\|_{TV} \leq \epsilon$, then for any bounded function f with $\sup |f| \leq M$, the difference in expectations satisfies $|\mathbb{E}_P[f] - \mathbb{E}_Q[f]| \leq \epsilon \cdot M$. In our context, the joint distributions of the original and surrogate blocks have total variation distance bounded by $B\beta(d)$ from Lemma 1, while the function being evaluated is the average loss $\frac{1}{B} \sum_{j=1}^B L_j$. Since our loss function is bounded by 1, each block average L_j is also bounded by 1, giving us $\sup_j |L_j| \leq 1$ and yielding the final bound of $B\beta(d)$. This bound establishes our mathematical bridge between dependent and independent learning, showing that when $B\beta(d)$ is small—achieved by setting d appropriately for exponentially decaying $\beta(d)$ —the error from our independence approximation becomes negligible. The inequality directly connects approximation quality to the mixing properties through $\beta(d)$, explicitly quantifying how temporal dependencies affect generalization.

This surrogate i.i.d. sequence is crucial because it allows us to apply the standard online-to-batch conversion result from Cesa-Bianchi and Lugosi Cesa-Bianchi & Lugosi (2006). Let \tilde{R}_B be the regret with respect to these surrogate variables, defined as $\tilde{R}_B = \sum_{j=1}^B \tilde{L}_j - \min_f \sum_{j=1}^B \ell(f, \tilde{Z}_j)$ where \tilde{Z}_j represents the surrogate data in block j .

Their theorem guarantees that for i.i.d. random variables, with probability at least $1 - \delta/2$:

$$\left| \frac{1}{B} \sum_{j=1}^B \tilde{L}_j - \mathbb{E}[\tilde{L}_1] \right| \leq \frac{\tilde{R}_B}{B} + \sqrt{\frac{\log(2/\delta)}{2B}}$$

This inequality represents the mathematical cornerstone of online-to-batch conversion, directly connecting the online learning regret to the generalization error. The term $\mathbb{E}[\tilde{L}_1]$ equals the true risk $\mathcal{L}(\tilde{f})$ (as established above using convexity), while $\frac{1}{B} \sum_{j=1}^B \tilde{L}_j$ approximates the empirical risk. Critically, the concentration term $\sqrt{\frac{\log(2/\delta)}{2B}}$ depends on the number of blocks B , not the original sequence length N . Since $B = \lfloor N/(d+1) \rfloor$ and we choose $d = \Theta(\log N)$ for exponential mixing, this concentration term scales as $\sqrt{\log(1/\delta) \log N / N}$ rather than the standard $\sqrt{\log(1/\delta)/N}$ rate for i.i.d. data. This $\sqrt{\log N}$ factor represents the price of converting dependent data into approximately independent blocks.

Step 3: Bounding error terms. In this final step, we connect our surrogate variables back to the original problem. We decompose the generalization error into manageable components and bound each separately.

First, we establish that $\mathbb{E}[\tilde{L}_1] = \mathcal{L}(\tilde{f})$, which relies crucially on convexity of the loss function. By construction, \tilde{L}_1 has the same marginal distribution as L_1 , so $\mathbb{E}[\tilde{L}_1] = \mathbb{E}[L_1]$. Now, $L_1 = \frac{1}{d+1} \sum_{t \in I_1} \ell(h_t, Z_t)$. By stationarity of the process, all time indices have identical marginal distributions, so:

$$\begin{aligned} \mathbb{E}[L_1] &= \mathbb{E} \left[\frac{1}{d+1} \sum_{t \in I_1} \ell(h_t, Z_t) \right] = \frac{1}{d+1} \sum_{t \in I_1} \mathbb{E}[\ell(h_t, Z_t)] \\ &= \frac{1}{d+1} \sum_{t \in I_1} \mathbb{E}[\ell(h_t, Z_1)] = \mathbb{E} \left[\frac{1}{d+1} \sum_{t \in I_1} \ell(h_t, Z_1) \right] \end{aligned}$$

Since the loss ℓ is convex in its first argument and $\tilde{f} = \frac{1}{N} \sum_{t=1}^N h_t$, by Jensen's inequality applied conditionally on Z_1 :

$$\ell(\tilde{f}, Z_1) \leq \frac{1}{N} \sum_{t=1}^N \ell(h_t, Z_1)$$

Taking expectations and using that the block I_1 is representative of the full sequence by stationarity, we obtain $\mathbb{E}[\ell(\bar{f}, Z_1)] = \mathbb{E}[L_1]$. Thus $\mathbb{E}[\tilde{L}_1] = \mathcal{L}(\bar{f})$, where the equality fundamentally depends on loss convexity.

The surrogate regret is bounded by the original regret scaled by the block size: $\tilde{R}_B \leq \frac{R_N}{d+1}$.

We can now decompose the generalization error:

$$\begin{aligned} |\mathcal{L}(\bar{f}) - \hat{\mathcal{L}}_N(\bar{f})| &\leq \left| \mathcal{L}(\bar{f}) - \frac{1}{B} \sum_{j=1}^B L_j \right| + \left| \frac{1}{B} \sum_{j=1}^B L_j - \hat{\mathcal{L}}_N(\bar{f}) \right| \\ &\leq B\beta(d) + \frac{R_N}{B(d+1)} + \sqrt{\frac{\log(2/\delta)}{2B}} + \frac{r}{N} \end{aligned}$$

The first term represents the approximation error from the coupling, the second comes from the online regret, the third is the concentration term for i.i.d. variables, and the fourth accounts for the remainder blocks.

Since $B = \lfloor N/(d+1) \rfloor$, we have $B(d+1) \leq N$ and $B \geq \frac{N}{d+1} - 1$. For exponential mixing with delay $d = \lceil \log N/c_0 \rceil$, we have $B \approx N/\log N$. The remainder term satisfies $r = N - B(d+1) < d+1$, so $\frac{r}{N} < \frac{d+1}{N}$.

The concentration term becomes $\sqrt{\frac{\log(2/\delta)}{2B}} \approx \sqrt{\frac{\log(1/\delta) \log N}{N}}$. With these bounds and using $N\beta(d) \geq B\beta(d)$, we obtain:

$$|\mathcal{L}(\bar{f}) - \hat{\mathcal{L}}_N(\bar{f})| \leq \frac{R_N}{N} + N\beta(d) + \sqrt{\frac{\log(1/\delta) \log N}{N}} + \frac{d+1}{N}$$

Absorbing the smaller terms, our final result is:

$$|\mathcal{L}(\bar{f}) - \hat{\mathcal{L}}_N(\bar{f})| = O\left(\frac{R_N}{N} + N\beta(d) + \sqrt{\frac{\log N}{N}}\right)$$

This bound has a clear interpretation: the generalization error is controlled by three terms—the average regret $\frac{R_N}{N}$ measuring optimization quality, the mixing term $N\beta(d)$ quantifying the effect of temporal dependencies, and a concentration term $\sqrt{\frac{\log N}{N}}$ that includes the $\sqrt{\log N}$ factor arising from block-level concentration with $B = O(N/\log N)$ blocks. When $\beta(d)$ decays exponentially with d , as in our Assumption 1, we can set $d = \Theta(\log N)$ to make the mixing term $N\beta(d) = O(1)$, effectively eliminating the impact of dependencies on the asymptotic convergence rate. This result shows that despite temporal dependencies, we can achieve generalization rates nearly identical to the i.i.d. case by properly spacing our observations and leveraging online learning techniques.

B.3 Proof of Lemma 2 (TCN Rademacher Complexity)

This lemma sets a non-vacuous bound on the capacity of TCNs to fit random noise—a key component in quantifying how architecture affects generalization. Intuitively, Rademacher complexity measures a hypothesis class’s ability to correlate with random patterns; lower complexity implies better generalization. Our goal is to show that despite their expressivity, TCNs with controlled architectural parameters maintain manageable complexity. We need to bound the Rademacher complexity of the class $\mathcal{F}_{D,p,R}$ consisting of TCNs with depth D , kernel size p , and weight norm bound R . The challenge lies in accounting for the convolutional structure and depth while avoiding exponential dependence on architectural parameters.

Step 1: Base layer analysis. We begin by analyzing a single-layer network as our foundation. For a single layer with input dimension n , the Rademacher complexity can be bounded using standard results for linear predictors with bounded norm. Since we constrain the $\ell_{2,1}$ norm by R , the Rademacher complexity of the base layer is bounded by:

$$\mathfrak{R}_m(\mathcal{F}_{1,p,R}) \leq R \sqrt{\frac{n}{m}}.$$

This bound encapsulates a fundamental statistical principle: complexity scales with the square root of input dimension n (reflecting the model’s capacity) and inversely with the square root of sample size m (reflecting the benefit of additional data). The parameter R acts as a multiplier—larger weight norms directly increase complexity and risk of overfitting.

Step 2: Layer-wise Lipschitz constants. To extend it to deeper networks, we examine how each layer transforms its inputs. Each convolutional layer with kernel size p and $\ell_{2,1}$ norm bounded by R is Lipschitz continuous with respect to its inputs. Lipschitz continuity quantifies, in essence, how much a layer can amplify small changes in its input—a critical property for understanding error propagation through neural networks.

Following the spectral analysis of convolutional operators by Sedghi et al. (2019), the Lipschitz constant can be bounded by $R\sqrt{p}$. This factor has an intuitive interpretation: \sqrt{p} appears because each output position depends on p consecutive input positions, creating potential for signal amplification proportional to the square root of the receptive field size. Meanwhile, the R factor reflects our weight constraint, which limits the sum of the Euclidean norms of the filters.

Step 3: Composition via contraction principle. A natural approach for deep networks is to compound the complexity layer by layer. For a composition of Lipschitz functions, a naive application of the chain rule would multiply the Lipschitz constants, giving a bound that grows exponentially with depth as $(R\sqrt{p})^D$. This would render the bound vacuous for even moderately deep networks—an issue that has historically plagued generalization theory for deep learning. To overcome this limitation, we leverage the vector contraction principle from empirical process theory Ledoux & Talagrand (2013) together with the Heinz-Khinchin smoothing techniques developed by Golowich et al. (2018). These advanced tools allow us to “smooth” the composition, avoiding exponential explosion with depth.

Let f_W denote a TCN in $\mathcal{F}_{D,p,R}$. We can view it as a composition $f_W = f_D \circ \dots \circ f_1$, where each f_i is a convolutional layer followed by a ReLU activation. By the vector contraction principle and the properties of ReLU activations (which are 1-Lipschitz and preserve the origin), we have:

$$\mathfrak{R}_m(\mathcal{F}_{D,p,R}) \leq 2\sqrt{D} \cdot (R\sqrt{p})^D \cdot \mathfrak{R}_m(\mathcal{F}_0)$$

where \mathcal{F}_0 is the class of identity functions on the input. However, this bound still contains the exponential term $(R\sqrt{p})^D$, which we need to improve further.

Step 4: Improved bound via Golowich et al. technique. The key insight from Golowich et al. Golowich et al. (2018) is that the *architectural depth factor* can be improved from exponential (c^D) to square-root (\sqrt{D}) dependence. However, this improvement applies specifically to how depth appears in the bound—it does *not* eliminate the product of layer-wise weight norms. More precisely, Golowich et al. achieve bounds of the form $(\prod_{\ell=1}^D M^{(\ell)}) \cdot \sqrt{D}$ rather than $(\prod_{\ell=1}^D M^{(\ell)})^D$ that would arise from naive Lipschitz composition. The product $\prod_{\ell=1}^D M^{(\ell)}$ necessarily remains because each layer’s contribution to the overall function depends on its weight magnitude.

For our hypothesis class $\mathcal{F}_{D,p,R}$, we constrain each layer to satisfy $\|W^{(\ell)}\|_{2,1} \leq M^{(\ell)}$ and define $R = \prod_{\ell=1}^D M^{(\ell)}$ as the product of these layer-wise bounds. The final Rademacher complexity bound will include R explicitly, representing the product of norms, multiplied by the improved \sqrt{D} depth factor.

Instead of analyzing layers sequentially (which compounds errors and creates exponential depth dependence in both the product and the depth), Golowich’s technique examines the network holistically. The key insight is that random noise does not simply accumulate as it propagates through a deep network—rather, cancellation effects occur between layers because random patterns rarely align consistently across all layers. These cancellation effects mean that the network’s capacity to fit random noise grows only with the square root of depth (\sqrt{D}) rather than exponentially (c^D). This improvement transforms our bounds from purely theoretical to practically meaningful—for a 9-layer network, complexity scales with 3 instead of $2^9 = 512$, making deep learning theory relevant to real-world architectures.

Technically, this approach involves analyzing the expected supremum of a Rademacher process indexed by the function class and applying martingale concentration inequalities that capture the dependencies between

layers. Adapting their approach to our convolutional setting while accounting for the specific structure of TCNs, we obtain:

$$\mathfrak{R}_m(\mathcal{F}_{D,p,R}) \leq 4R \sqrt{\frac{D p n \log(2m)}{m}}$$

where $R = \prod_{\ell=1}^D M^{(\ell)}$ is the product of layer-wise weight norm bounds. This bound achieves \sqrt{D} scaling for the architectural depth factor (improving from exponential R^D that would arise from naive composition) while necessarily retaining the product of layer norms R .

The factor $\log(2m)$ appears from covering number arguments used in the proof and grows very slowly with sample size. The final bound reveals how each architectural parameter contributes to model complexity:

- The linear dependence on R shows that weight norm directly scales complexity
- The square root dependence on depth D demonstrates that deeper networks increase complexity much more slowly than an exponential relationship would suggest
- The square root dependence on kernel size p indicates that larger receptive fields increase complexity, as each output draws information from more inputs
- The factor $\sqrt{n/m}$ shows the relationship between input (n) dimension and sample size (m)

This result is significant because it proves that even deep TCNs can generalize well with sufficient data, overcoming the pessimistic predictions of naive bounds. The explicit dependence on architectural parameters provides practical guidance for model design: doubling the depth increases complexity by only about 40%, while doubling the kernel size increases complexity by about 40% as well (using the square root). This bound plays a central role in our main generalization theorem by quantifying exactly how architectural choices affect learning from dependent data. It ensures that our bounds remain meaningful even for the deep models used in contemporary applications.

B.4 Proof of Theorem 1 (Architecture-Aware Generalization)

This theorem represents the culmination of our work, combining all previous results to provide explicit, architecture-aware bounds for temporal models trained on dependent data. We will show how the optimal choice of delay parameter, combined with our Rademacher complexity results, leads to practical generalization guarantees that explicitly depend on network architecture.

Step 1: Optimal delay parameter selection. The first key insight is how to choose the delay parameter d to effectively balance between reducing dependencies and maintaining sufficient training data. Under Assumption 1, the β -mixing coefficient satisfies $\beta(k) \leq C_0 e^{-c_0 k}$ for constants $C_0, c_0 > 0$.

We set $d = \lceil \log N / c_0 \rceil$, which has a clear intuitive meaning: we choose a delay that grows logarithmically with sequence length, with the specific rate determined by the mixing rate c_0 of the underlying process. This choice allows us to show:

$$\begin{aligned} N\beta(d) &\leq N \cdot C_0 e^{-c_0 d} \\ &\leq N \cdot C_0 e^{-c_0 \lceil \log N / c_0 \rceil} \\ &\leq N \cdot C_0 e^{-\log N} \\ &= C_0 \end{aligned}$$

This calculation reveals that with our logarithmic choice of delay, the mixing-dependent term $N\beta(d)$ becomes a constant C_0 independent of the sample size. This effectively eliminates the impact of temporal dependencies on the asymptotic learning rate, allowing our dependent-data bound to match the structure of classical i.i.d. bounds.

Step 2: Regret bound via Rademacher complexity. Next, we leverage our Rademacher complexity bound from Lemma 2 to bound the regret term in Proposition 1. For our hypothesis class $\mathcal{F}_{D,p,R}$ of TCNs with depth D , kernel size p , and weight norm bound R , we have:

$$\mathfrak{R}_m(\mathcal{F}_{D,p,R}) \leq 4R \sqrt{\frac{D p n \log(2m)}{m}}$$

For online convex optimization with mirror descent using an ℓ_2 -regularizer and step size $\eta_t = \sqrt{\log N/t}$, standard results (e.g., Shalev-Shwartz 2012) bound the regret in terms of Rademacher complexity. The connection $R_N \leq N \cdot \mathfrak{R}_N(\mathcal{F})$ follows from online-to-batch conversion under convex losses. Note that this step again requires loss convexity—without it, we cannot bound regret via Rademacher complexity in this manner. Applying this with our Rademacher bound:

$$\begin{aligned} R_N &\leq 2N \mathfrak{R}_N(\mathcal{F}_{D,p,R}) \\ &\leq 8NR \sqrt{\frac{D p n \log(2N)}{N}} \\ &= 8R \sqrt{D p n N \log(2N)} \end{aligned}$$

where $R = \prod_{\ell=1}^D M^{(\ell)}$ is the product of layer-wise norms.

This gives us a bound on the total regret R_N . For the delayed-feedback approach, we need the per-sample regret:

$$\begin{aligned} \frac{R_N}{N} &\leq 8R \sqrt{\frac{D p n \log(2N)}{N}} \\ &= O\left(R \sqrt{\frac{D p n \log N}{N}}\right) \end{aligned}$$

This per-sample regret bound explicitly shows how model complexity (through D , p , and R) affects learning performance. The bound decreases with sample size N at rate $1/\sqrt{N}$, the optimal rate for non-parametric learning, while increasing with model complexity parameters in an interpretable way.

Step 3: Combining the bounds. We now apply Proposition 1 with our bounds for $N\beta(d)$ and R_N/N . For any $f \in \mathcal{F}_{D,p,R}$ produced by the delayed-feedback learner with convex Lipschitz loss, with probability at least $1 - \delta$:

$$\begin{aligned} |\mathcal{L}(f) - \widehat{\mathcal{L}}_N(f)| &\leq \frac{R_N}{N} + N\beta(d) + \sqrt{\frac{\log(1/\delta) \log N}{N}} \\ &\leq C_1 R \sqrt{\frac{D p n \log N}{N}} + C_0 + \sqrt{\frac{\log(1/\delta) \log N}{N}} \end{aligned}$$

where $R = \prod_{\ell=1}^D M^{(\ell)}$ is the product of layer-wise weight norms, $C_1 = 8$ from our regret analysis, and C_0 is the constant from exponential β -mixing. The concentration term includes $\sqrt{\log N}$ due to block-level analysis with $B = O(N/\log N)$ blocks.

Structure of the Corrected Bound. This bound explicitly includes:

- The product of layer norms $R = \prod_{\ell=1}^D M^{(\ell)}$ in the complexity term

- The improved \sqrt{D} architectural depth factor (not exponential c^D)
- A concentration term with $\sqrt{\log N}$ factor from block-level analysis
- A constant C_0 mixing term when $d = O(\log N)$ is chosen optimally

Under global norm constraints where $R \leq R_0$ for fixed constant R_0 , the dominant term becomes $O(\sqrt{D \log N}/N)$. This scaling suggests that doubling depth twice (from D to $4D$) gives factor $\sqrt{4} = 2$ in the complexity term, requiring approximately $4\times$ more data to maintain the same bound—though this is a worst-case prediction, and our experiments (Section 5.2) show architectures well-matched to temporal structure often achieve better sample efficiency.

This final bound has a clear interpretation: the generalization gap for TCNs trained on β -mixing data is controlled by four terms; 1. A complexity term $R\sqrt{\frac{D p n \log N}{N}}$ that explicitly shows how architectural parameters (D, p) and the product of layer norms (R) affect generalization; 2. A constant mixing term C_0 that captures the irreducible impact of temporal dependencies, and 3. A standard concentration term $\sqrt{\frac{\log(1/\delta)}{N}}$ reflecting the confidence parameter. The bound demonstrates that with sufficient data, even complex temporal models can generalize well on dependent data. Moreover, it provides practical guidance for architecture selection by quantifying exactly how different design choices impact generalization. This completes the proof of Theorem 1.