
Convergent Stochastic Training of Attention and Understanding LoRA

Anonymous Authors¹

Abstract

Transformers have revolutionized machine learning and deploying attention layers in the model is increasingly standard across a myriad of applications. Further, for large models, it is common to implement Low Rank Adaptation (LoRA), whereby a factorized parameterization of them is trained, to achieve a surprisingly beneficial accuracy-size trade-off. In this work, via a unified framework we rigorously establish trainability of such models under stochastic methods. We prove that for any mild regularization, the empirical regression loss on a attention layer and LoRA on a shallow neural net, both induce Poincaré inequality for the corresponding Gibbs’s measure. Then it follows via invoking recent results that a certain SDE, which mimics the SGD, minimizes the corresponding losses. In both the cases, our first-of-its-kind results of trainability on attention and nets, do not rely on any assumptions on the data or the size of the architecture.

1. Introduction

The remarkable empirical success of attention mechanisms have fundamentally reshaped modern machine learning, most prominently through the transformer architecture – which is the backbone of Large Language Models (LLMs) (Radford et al., 2018). By allowing models to dynamically weight interactions between “tokens”/data fragments, attention layers enable modeling of complicated distributions. Central to the attention mechanism are the query and key matrices, \mathbf{W}_Q and \mathbf{W}_K which occur in the model as the product $\mathbf{W}_Q \mathbf{W}_K$ and one such pair are trained in each “attention head” of which there are many in each of the many attention layers in any commonly used transformer. Despite their central role in the practice of modern AI, our theoretical understanding remains limited of how these matrices evolve during successful training.

While modern attention mechanisms are now the standard,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

they originally emerged to address the bottleneck of fixed-length representations in early sequence modeling. Early sequence-to-sequence models were limited by fixed-length representations. (Bahdanau et al., 2016) introduced attention to dynamically aggregate encoder states, allowing the decoder to focus on relevant input positions. Building on this, (Luong et al., 2015) proposed alternative attention variants, including global attention over all positions and local attention over a subset.

These developments ultimately led to the transformer architecture (Vaswani et al., 2017), which removes recurrence entirely and instead models sequence interactions purely through stacked self-attention layers, where given an input $\mathbf{X} \in \mathbb{R}^{t \times d}$ each layer computes $\text{RowSoftMax}_\beta \left(\frac{\mathbf{X} \mathbf{W}_Q \mathbf{W}_K^T \mathbf{X}^T}{\sqrt{d}} \right) \mathbf{X} \mathbf{W}_V$, where $\mathbf{W}_Q, \mathbf{W}_K$ and \mathbf{W}_V are learned weight matrices where the RowSoftMax operator is defined as $[\text{RowSoftMax}_\beta(\mathbf{M})]_{ij} := \frac{e^{\beta M_{ij}}}{\sum_{k=1}^t e^{\beta M_{ik}}}$.

The transformer architecture (Vaswani et al., 2017) has become the dominant paradigm for sequence modeling, replacing recurrence with stacked self-attention and feed-forward layers. Numerous variants have since been proposed, namely Sparse Transformer (Child et al., 2019), Longformer (Beltagy et al., 2020), Linformer (Wang et al., 2020), Transformer-XL (Dai et al., 2019), ALBERT (Lan et al., 2020), and Vision Transformer (Dosovitskiy et al., 2021), among many others. A foundational reason explaining these successes is given in works like (Yun et al., 2020) that have show that transformers are universal approximators of sequence-to-sequence matrix functions.

The Rising Importance of Doing Regression on Transformers The success of transformers has led to their widespread adoption in scientific machine learning, where many tasks can be naturally formulated as regression problems over high-dimensional discretizations of function spaces. In particular, applications in fluid dynamics and weather prediction — ranging from operator learning for PDEs to data-driven forecasting — often reduce to learning mappings between input and output fields by doing regression using a attention-based architecture, FourCastNet (Pathak et al., 2022), GraphCast (Lam et al., 2023), Pangu-Weather (Bi et al., 2023), Poseidon (Herde et al., 2024) and GenCFD (Molinario et al., 2024). This perspec-

055 tive motivates studying attention not only as a representation
 056 mechanism, but as a regression operator whose properties
 057 govern generalization and efficiency in continuous domains.

058 **Q1:** *With no assumptions on data or architecture, can a*
 059 *attention layer be trained by a stochastic algorithm?*
 060

061 On the other hand, the rise of large-scale pretraining has
 062 motivated parameter-efficient fine-tuning methods such as
 063 Low-Rank Adaptation (LoRA), which constrains updates
 064 to lie in low-dimensional subspaces while preserving the
 065 pretrained backbone. More precisely, the idea is to freeze
 066 the pre-trained weight matrix $\mathbf{W}_{\text{pre}} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ and only train
 067 a small update. This update is represented as the product of
 068 two small matrices, \mathbf{A} and \mathbf{B} . The modified weights \mathbf{W}' are
 069 parameterized as, $\mathbf{W}' = \mathbf{W}_{\text{pre}} + \frac{\alpha}{r} \mathbf{B} \mathbf{A}$ where $\mathbf{A} \in \mathbb{R}^{r \times d_{\text{in}}}$
 070 and $\mathbf{B} \in \mathbb{R}^{d_{\text{out}} \times r}$ are the trainable factors and α being a
 071 scaling constant. By training only \mathbf{A} and \mathbf{B} , we significantly
 072 reduce the number of parameters to update.
 073

074 As showed in (Shuttleworth et al., 2025), LoRA works best
 075 when applied to all weight matrices and small-to-medium
 076 datasets, and its optimal learning rate is largely independent
 077 of rank due to the $\frac{1}{r}$ scaling. Overall, LoRA performs
 078 similar to full-finetuning in the “low-regret regime” making
 079 it a parameter-efficient alternative for post-training adap-
 080 tation. However, the product $\mathbf{B} \mathbf{A}$ introduces a specific
 081 scaling problem : we can multiply \mathbf{A} by a constant and
 082 divide \mathbf{B} by the same constant without changing the final
 083 output. This redundancy makes the training landscape “flat”
 084 in certain directions, which is seemingly a natural obstacle
 085 for gradient based algorithms to succeed. And yet LoRA
 086 has proven strikingly effective in practice, dramatically re-
 087 ducing memory and computational overhead while retaining
 088 performance. Thus we posit that the optimization dynamics
 089 of neural models under LoRA have remained unclear from
 090 a theoretical standpoint.

091 **Q2:** *With no assumptions on data or width, can a*
 092 *implementation of LoRA on a net be trained by a stochastic*
 093 *algorithm?*
 094

095 In this work, *firstly* we make progress towards uncover-
 096 ing hitherto unknown mathematical properties about the
 097 attention map and thus uncover a first-of-its-kind provably
 098 convergent training mechanism for the query and key matri-
 099 ces. *Secondly*, this work also initiates a theoretical study of
 100 training dynamics for standard neural networks under LoRA
 101 parameterization — which evidently shares a mathematical
 102 similarity of training a factorized weight parameterization
 103 as while training the key and query matrices in the attention,
 104 as introduced above.

105 Whether for training the key and query matrices of an at-
 106 tention head or for training standard nets under LoRA, we
 107 consider training through stochastic differential equations
 108 (SDEs), capturing the continuous-time limit of stochastic
 109

gradient methods commonly used in practice. *In both cases*
— for any number of parameters and for any data — we
prove convergence of a risk function for the models under a
SDE flow.

1.1. Summary of Results

In this work, we establish that for both attention and depth-2
 neural networks under LoRA, a mildly regularized regres-
 sion loss (say \tilde{V}) is a Villani function (which will be pre-
 cisely defined in the next section) — which in turn implies
 that the corresponding Gibbs’ measure ($\sim e^{-\gamma \tilde{V}}$) satisfies
 the Poincare inequality. Then invoking recent results on
 isoperimetry based Stochastic Differential Equation (SDE)
 convergence (Shi et al., 2023) we can establish convergence
 in both settings for the following continuous-time stochastic
 gradient dynamics for the weights \mathbf{T} , given by the SDE,

$$d\mathbf{T}_t = -\nabla \tilde{V}(\mathbf{T}_t) dt + \sqrt{s} d\mathbf{B}_t, \quad (1)$$

where $\tilde{V}(\mathbf{T})$ denotes the regularized loss function, $s > 0$
 is a temperature parameter and $(\mathbf{B}_t)_{t \geq 0}$ is the Brownian
 motion.¹ We provide an informal restatement of our main
 results below.

Theorem 1.1 (Informal Statement of Provable Learning for
 Attention-Based Regression). *Consider a single attention*
layer with key and query matrices \mathbf{W}_K and \mathbf{W}_Q , being
trained using the ℓ_2 -loss function with either a logarithmi-
cally amplified 2-norm regularization or a super-quadratic
polynomial regularization. Then, for any arbitrarily low
regularization, for any data and size of architecture, the loss
function satisfies the Villani condition.

As a consequence, for any $\varepsilon > 0$, there exists an appropriate
step size s such that the SDE for $\mathbf{T} = (\mathbf{W}_Q, \mathbf{W}_K)$ con-
verges, in expectation, in $\mathcal{O}(\log \frac{1}{\varepsilon})$ to within ε of the global
minimum of the training loss.

The above informal restatement combines Theorem 3.1 with
 the relevant part of Theorem 3.3.

Theorem 1.2 (Informal Statement of Provable Learning
 for Depth-2 Neural Net Based Regression under LoRA).
Consider a depth-2 neural network with weight matrix \mathbf{W}
factorized as $\mathbf{U}\mathbf{V}$, being trained using the ℓ_2 -loss function
with either a logarithmically amplified 2-norm regulariza-
tion or a super-quadratic polynomial regularization. Then,
for any arbitrarily low regularization, for any data and
size of architecture, the loss function satisfies the Villani
condition.

As a consequence, for any $\varepsilon > 0$, there exists an appropriate
step size s such that the SDE for $\mathbf{T} = (\mathbf{U}, \mathbf{V})$ converges, in
expectation, in $\mathcal{O}(\log \frac{1}{\varepsilon})$ to within ε of the global minimum
of the training loss.

¹The subscript t denotes continuous time.

The above informal restatement combines Theorem 3.2 with the relevant part of Theorem 3.3.

Remark 1.3. As shown in (Kumar et al., 2025), Theorems 1.1 and 1.2 can be extended to establish convergence to the global minima of the population risk, defined as $\mathbb{E}_{\mathcal{S}_n}[V_{\mathcal{S}_n}(\mathbf{T})]$, where $V_{\mathcal{S}_n}(\mathbf{T})$ denotes the loss evaluated on the dataset $\mathcal{S}_n = (\mathbf{x}_i, y_i)_{i=1}^n$, under the Langevin Monte Carlo algorithm.

1.2. Literature Review

A theoretical analysis of transformers in the infinite-width limit by deriving their Neural Network Gaussian Process (NNGP) and Neural Tangent Kernel (NTK) equivalents was provided in (Hron et al., 2020). They showed that self-attention layers admit well-defined kernel limit and in this regime, gradient descent training of a Transformer is equivalent to kernel regression with the corresponding NTK.

A global convergence framework for transformers by analyzing training dynamics in the mean-field limit was established in (Gao et al., 2024). By treating model width and depth as approaching infinity, the authors demonstrate that discrete gradient descent converges to a Wasserstein gradient flow on the distribution of parameters. Albeit the use of an infinite-width limit it strictly requires the inclusion of a weight decay parameter $\lambda > 0$, for the convergence proof to work.

A formal proof that transformers can converge to the functional behavior of near-optimal Reinforcement Learning (RL) algorithms through the optimization of a log-likelihood objective was given by (Lin et al., 2023). By viewing the attention mechanism as an iterative optimizer, they demonstrate that supervised pre-training on offline trajectories allows the model to implement algorithms like LinUCB and Thompson Sampling directly. However, these convergence guarantees rely on non-standard architectures, most notably the use of ReLU-based attention to facilitate exact linear algebraic operations. Furthermore, the authors assume model realizability, implying that the Transformer’s capacity must be sufficient to encapsulate the expert’s decision-making logic.

A multi-layer transformer trained on n -gram data, where each token depends on the preceding n tokens, is analyzed in (Chen et al., 2024), and it is shown that gradient flow converges to a model exhibiting induction head behavior. Specifically, induction here refers to the phenomenon whereby, if a token at position i matches a previous occurrence at position j , the model attends to the token following position j to predict the next token at position i . Their results provide a rigorous characterization of how attention layers, feed-forward networks, and normalization interact to learn features from context. This work advances prior studies that focused on linear or single-layer models by handling

richer architectures and more realistic data distributions.

Fundamental algorithmic limits of Multi-Head Attention (MHA) restricted to a discrete Boolean input distribution $X \in \{\pm 1\}^{k \times d}$ was established in (Chen & Li, 2025). Theorem 1.2 in (Chen & Li, 2025) establishes that, under a non-degeneracy condition on the attention and projection matrices and realizability assumption of the samples, there exists an algorithm that for m -headed attention estimates the parameters in $(kd)^{O(m^3)}$ time, using $(kd)^{\Theta(m)}$ samples, and achieves predictions that are $(kd)^{-\Omega(m)}$ close to the true values in expectation. This result identifies the number of heads m as the dominant factor in computational scaling.

1.2.1. REVIEW OF EXISTING ATTEMPTS AT PROVABLE TRAINING OF LORA ON NEURAL NETWORKS — WITH WEIGHT REGULARIZATION

LoRA was first introduced by (Hu et al., 2022) where it was asserted that the weight updates for task-specific adaptation in attention based models reside in a manifold of low intrinsic dimension (Aghajanyan et al., 2021). By reparameterizing the update matrix $\Delta \mathbf{W}$ as the product of two low-rank matrices \mathbf{A} and \mathbf{B} , (Hu et al., 2022) demonstrated that optimization can converge to high-performance solutions with significantly fewer trainable parameters. Crucially, their initialization strategy — setting one matrix to zero — ensures a stable starting point at the pre-trained state, effectively bridging the gap between training efficiency and the convergence stability typically observed in full fine-tuning.

A rigorous analysis of LoRA on neural networks in the generic non-linear regime was given by (Kim et al., 2025), establishing a formal dichotomy between global convergence and parameter divergence. By characterizing the optimization as governed by a global Restricted Strong Convexity, they prove that the combination of zero-initialization and weight decay induces an implicit bias toward a low-rank global minimum. Crucially, they demonstrate that while the non-linear landscape may harbor spurious local minima, these points are spectrally isolated in high-rank regions and do not intersect with the stable optimization trajectory.

In the Neural Tangent Kernel (NTK) framework LoRA was analyzed by (Jang et al., 2024), where the neural network’s optimization can be treated as a linearized system. In this regime, the training dynamics are governed by a quadratic objective subject to a low-rank structural constraint. Alongside a non-standard regularizer — derived from the Rademacher complexity of the low-rank bottleneck — the authors prove that the non-convex BA reparameterization does not introduce spurious local minima. (Jang et al., 2024) further demonstrates, if the rank is above a certain threshold, gradient-based methods like stochastic gradient descent (SGD) converge to a low-rank global minimizer.

Limitations of LoRA It was demonstrated by (Shuttleworth et al., 2025) that the perceived equivalence between LoRA and full fine-tuning is an “illusion” maintained by surface-level metrics. Despite its parameter efficiency, LoRA can lead to catastrophic forgetting, thereby degrading performance in settings that require continual learning.

1.2.2. ORGANIZATION

Section 2 introduces the analytic framework, including the underlying conditions, neural architecture, loss functions, and the assumptions required for the subsequent sections. Section 3 presents the main results of the paper, namely Theorem 3.1, Theorem 3.2, and Theorem 3.3. The proof of Theorem 3.1, along with the requisite auxiliary lemmas, is provided in Section A, while the detailed proofs of these lemmas are deferred to Section B. Section 4 presents experimental results on solving the 2D Darcy flow problem using the regularized loss functions introduced in Section 2. The proof of Theorem 3.2 is contained in Appendix C, with the corresponding supporting lemmas established in Appendix D. Lastly, Appendix E presents the proof of Theorem 3.3.

2. Mathematical Setup

In this section, we define the analytic conditions, architectures, and loss functions, alongside the core assumptions that underpin our subsequent analysis.

The Villani condition was introduced in (Villani, 2009) to guarantee that, when a function satisfies this condition, the associated Gibbs measure satisfies the Poincaré inequality. We recall that a distribution π is said to satisfy the Poincaré inequality for some constant C_{PI} , if for all smooth functions $h : \mathbb{R}^d \rightarrow \mathbb{R}$, $\text{Var}_\pi(h) \leq C_{PI} \mathbb{E}_\pi[\|\nabla h\|^2]$. (Shi et al., 2023) leverage the Poincaré inequality induced by the Villani condition to establish convergence results for certain stochastic differential equations (SDEs). We now proceed to formally define the corresponding analytic conditions.

Definition 2.1 (Confining Condition). A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be confining if it satisfies the following conditions, **(1.)** $f \in C^\infty$, **(2.)** $\lim_{\|\mathbf{x}\| \rightarrow +\infty} f(\mathbf{x}) = +\infty$, and **(3.)** $\int_{\mathbb{R}^d} e^{-\frac{2f(\mathbf{x})}{s}} d\mathbf{x} < \infty \forall s > 0$.

Definition 2.2 (Villani Condition). A confining function f is said to satisfy the Villani condition if for all $s > 0$, $\frac{\|\nabla f(\mathbf{x})\|^2}{s} - \Delta f(\mathbf{x}) \rightarrow \infty$ as $\|\mathbf{x}\| \rightarrow \infty$:

In recent works (Gopalani et al., 2024; Gopalani & Mukherjee, 2025), it was shown that the Villani condition holds for depth-2 neural networks of arbitrary width and for both squared and logistic losses. Building on this, (Kumar et al., 2025) proved that Langevin Monte Carlo consequently achieves population risk minimization. In contrast to (Gopalani et al., 2024; Gopalani & Mukherjee, 2025;

Kumar et al., 2025), our results do not require any lower bound on the regularization parameter.

Next we formally define the attention model and its regression loss that we choose to train.

Definition 2.3 (Attention Layer).

$$\begin{aligned} \mathbb{R}^{t \times d} \ni \mathbf{X} &\mapsto \text{Attention}(\mathbf{X}) \\ &:= \text{RowSoftMax}_\beta \left(\frac{\mathbf{X} \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}^\top}{\sqrt{d}} \right) \mathbf{X} \mathbf{W}_v \in \mathbb{R}^{t \times d} \end{aligned} \quad (2)$$

where, $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times r}$ and $\mathbf{W}_v \in \mathbb{R}^{d \times d}$

and for any $\mathbf{M} \in \mathbb{R}^{t \times t}$,

$$[\text{RowSoftMax}_\beta(\mathbf{M})]_{ij} = \frac{e^{\beta M_{ij}}}{\sum_{k=1}^t e^{\beta M_{ik}}}$$

A notable application of this framework is the Vision Transformer (ViT) (Dosovitskiy et al., 2021). In ViT, the input image is first divided into non-overlapping patches, commonly of size 16×16 pixels. Each patch is flattened and projected into a d -dimensional embedding, producing a sequence of t token embeddings that serve as input to the attention layers described above. This allows the model to capture long-range dependencies across the image while leveraging the same attention mechanism as in general Transformer architectures. The row-wise softmax scaling parameter β is typically set to 1.

We train the above model we consider two forms of factor-regularized potentials/loss functions for it : a non-polynomial/logarithmically amplified 2-norm regularization and a polynomial regularization with exponent $2 + \epsilon$. The following definitions formalize the corresponding regularized potentials used in our analysis.

Definition 2.4 (Mean Square Loss on a Attention Layer with Non-Polynomial Factor-Regularization). We define the potential $V_{\text{ATT}}(\mathbf{T})$ in the factor space \mathbb{R}^D for $\mathbf{T} = (\mathbf{W}_Q, \mathbf{W}_K)$ and $D = 2dr$ as,

$$\begin{aligned} \tilde{V}_{\text{ATT}}(\mathbf{T}) &:= \hat{R}_A(\mathbf{T}) \\ &+ \frac{\lambda}{2} (\|\mathbf{W}_Q\|_F^2 + \|\mathbf{W}_K\|_F^2) \log(1 + \|\mathbf{W}_Q\|_F^2 + \|\mathbf{W}_K\|_F^2) \end{aligned} \quad (3)$$

where $\hat{R}_A(\mathbf{T}) := \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{W}_Q, \mathbf{W}_K)$ and $\ell_i(\mathbf{W}_Q, \mathbf{W}_K) := \frac{1}{2} \left\| \mathbf{Y}_i - \text{RowSoftMax}_\beta \left(\frac{\mathbf{X}_i \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}_i^\top}{\sqrt{d}} \right) \mathbf{X}_i \mathbf{W}_v \right\|_F^2$ corresponding to a choice of training data as, $\{(\mathbf{X}_i, \mathbf{Y}_i) \in \mathbb{R}^{t \times d} \times \mathbb{R}^{t \times d} \mid i = 1, \dots, n\}$.

Definition 2.5 (Mean Square Loss on a Attention Layer with Polynomial Factor-Regularization). We define the potential $V_{\epsilon, \text{ATT}}(\mathbf{T})$ in the factor space \mathbb{R}^D for $\mathbf{T} =$

($\mathbf{W}_Q, \mathbf{W}_K$) and $D = 2dr$ as,

$$V_{\epsilon, \text{ATT}}(\mathbf{T}) := \hat{R}_A(\mathbf{T}) + \frac{\lambda}{2} (\|\mathbf{W}_Q\|_F^{2+\epsilon} + \|\mathbf{W}_K\|_F^{2+\epsilon}) \quad (4)$$

where $\hat{R}_A(\mathbf{T})$ as given in Definition 2.4.

For the case of shallow neural networks, we assume a training in a space of weight matrices with a rank bound i.e we assume the trainable weight to be factorizable as, $\mathbf{W} = \mathbf{U}\mathbf{V}^\top$, with $\mathbf{U} \in \mathbb{R}^{p \times r}$ and $\mathbf{V} \in \mathbb{R}^{d \times r}$. This factorization implements low-rank adaptation (LoRA) approach for depth-2 nets — and we recall that this has been previously studied in (Jang et al., 2024; Kim et al., 2025) as a theoretical sandbox for the LoRA technique. As in the attention-based setting, regularization of the factor matrices is introduced to ensure well-behaved potentials in the weight space. Similar to the attention training setup above, we define two types of factor-regularized loss functions as follows,

Definition 2.6 (Rank-Restricted Mean Square Loss on Shallow Nets with Non-Polynomial Factor-Regularization). We define the potential $V(\mathbf{T})$ in the factor space \mathbb{R}^D for $\mathbf{T} = (\mathbf{U}, \mathbf{V})$ and $D = (p+d)r$ as,

$$V(\mathbf{T}) := \mathcal{L}(\mathbf{T}) + \frac{\lambda}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) \log(1 + \|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) \quad (5)$$

where $\mathcal{L}(\mathbf{T}) := \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{U}\mathbf{V}^\top)$ and $\ell_i(\mathbf{W}) := \frac{1}{2} (y_i - \mathbf{a}^\top \sigma(\mathbf{W}\mathbf{x}_i))^2$.

Definition 2.7 (Rank-Restricted Mean Square Loss on Shallow Nets with Polynomial Factor-Regularization). We define the potential $V_\epsilon(\mathbf{T})$ in the factor space \mathbb{R}^D for $\mathbf{T} = (\mathbf{U}, \mathbf{V})$ and $D = (p+d)r$ as,

$$V_\epsilon(\mathbf{T}) = \mathcal{L}(\mathbf{T}) + \frac{\lambda}{2} (\|\mathbf{U}\|_F^{2+\epsilon} + \|\mathbf{V}\|_F^{2+\epsilon}) \quad (6)$$

where $\mathcal{L}(\mathbf{T}) := \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{U}\mathbf{V}^\top)$ and $\ell_i(\mathbf{W}) := \frac{1}{2} (y_i - \mathbf{a}^\top \sigma(\mathbf{W}\mathbf{x}_i))^2$.

To establish convergence results for the factorized, regularized losses defined above, we impose a few standard assumptions on the network and define certain properties of the training data.

Definition 2.8 (Training Data Bounds). The training data is bounded as follows:

1. For Attention : Each training example $(\mathbf{X}_i, \mathbf{Y}_i) \in \mathbb{R}^{t \times d} \times \mathbb{R}^{t \times d}$ satisfies $\|\mathbf{X}_i\|_F \leq B_x, \|\mathbf{Y}_i\|_F \leq B_y, i = 1, \dots, n$.

2. For Neural Network : Each training example $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$ satisfies $\|\mathbf{x}_i\|_2 \leq B_x, |y_i| \leq B_y, i = 1, \dots, n$.

Depending on the architecture, the bounds B_x and B_y are interpreted according to the corresponding definitions given above.

Definition 2.9 (Attention-Specific Bound). For regression tasks on attention-based models, the weight matrix \mathbf{W}_v is also bounded as $\|\mathbf{W}_v\|_F \leq B_w$.

Assumption 2.10 (Activation Function Bounds). We assume that σ', σ' and σ'' are bounded by $\sup_{t \in \mathbb{R}} |\sigma(t)| = B_\sigma < \infty, \sup_{t \in \mathbb{R}} |\sigma'(t)| = B_{\sigma'} < \infty$ and $\sup_{t \in \mathbb{R}} |\sigma''(t)| = B_{\sigma''} < \infty$, respectively.

We can then characterize the global convergence of the SDE (1), motivated by (Shi et al., 2023) in their analysis of SGD on non-convex landscapes, for our regularized loss function for depth-2 nets under LoRA constraints and for the attention-based model. The SDE is modeled over \mathbf{T} . The corresponding invariant Gibbs measure $\mu_s(d\mathbf{T})$ is defined as

$$\mu_s(d\mathbf{T}) := \frac{1}{Z_s} \exp\left(-\frac{2}{s} \tilde{V}(\mathbf{T})\right) d\mathbf{T}, \quad (7)$$

where $\tilde{V}(\mathbf{T})$ represents any of the factor-regularized potentials defined above, $s > 0$ acts as the temperature parameter (proportional to the learning rate), and Z_s is the normalization constant.

3. Main Results

Given the formal setup in the previous section, we first state our key result showing that the regression loss functions associated with the softmax-attention layer, defined in Definitions 2.4 and 2.5, satisfy the Villani condition.

Theorem 3.1 (Attention-Based Regression Loss is a Villani Function). Consider the regularized loss functions associated with the attention layer, $\tilde{V}_{\text{ATT}}(\mathbf{T})$ and $V_{\epsilon, \text{ATT}}(\mathbf{T})$, as defined in Definitions 2.4 and 2.5, respectively. Then, for any $\lambda, \epsilon > 0$, both $\tilde{V}_{\text{ATT}}(\mathbf{T})$ and $V_{\epsilon, \text{ATT}}(\mathbf{T})$, evaluated on the training data defined in Definition 2.8, satisfy the Villani condition (Definition 2.2).

We next present our second key result, establishing that the regression loss functions for a depth-2 neural network with LoRA, as defined in Definitions 2.6 and 2.7, satisfy the Villani condition.

Theorem 3.2 (Depth-2 Neural Net Based Regression Loss under LoRA is a Villani Function). Suppose that Assumptions 2.10 holds for the activation function σ . Consider the loss functions $\tilde{V}(\mathbf{T})$ and $V_\epsilon(\mathbf{T})$ for a depth-2 neural network with activation σ , as defined in Definitions 2.6 and 2.7, respectively. Then, for any $\lambda, \epsilon > 0$, both $\tilde{V}(\mathbf{T})$ and $V_\epsilon(\mathbf{T})$, evaluated on the training data defined in Definition 2.8, satisfy the Villani condition (Definition 2.2).

By Theorems 3.1 and 3.2, all considered loss functions with their associated neural architectures satisfy the Villani condition. Consequently, we may invoke Theorem 1 of (Shi et al., 2023) to obtain the following convergence result for the SDE (1).

Theorem 3.3 (Convergence of SDE for Depth-2 Neural Net Based Regression under LoRA and Attention-Based Regression). *Suppose that Assumption 2.10 holds for the activation σ . Let $\tilde{L}^{(k)}(\mathbf{T})$ denote any of the four regularized potentials defined in Definition 2.4 (\tilde{V}_{ATT}), Definition 2.5 ($V_{\epsilon, \text{ATT}}$), Definition 2.6 (V), and Definition 2.7 (V_ϵ), where $k \in \{1, 2, 3, 4\}$ indexes the specific model and regularization choice. Suppose the initial probability density is $p_0 \in L^2((\mu_s^{(k)})^{-1})$ of the SDE (1) where $\mu_s^{(k)}$ is the corresponding Gibbs measure (7). For each $\tilde{L}^{(k)}(\mathbf{T})$ satisfying the Villani conditions, there exists a positive $\lambda_s^{(k)} > 0$ and a constant $D^{(k)}(s, p_0)$ such that:*

$$\mathbb{E}[\tilde{L}^{(k)}(\mathbf{T}_t)] - \tilde{L}^{(k)*} \leq \epsilon^{(k)}(s) + D^{(k)}(s, p_0)e^{-\lambda_s^{(k)}t}, \quad (8)$$

where $\tilde{L}^{(k)*} = \inf \tilde{L}^{(k)}(\mathbf{T})$ is the global minimum of the respective loss, and $\epsilon^{(k)}(s) = \mathbb{E}_{\mu_s^{(k)}}[\tilde{L}^{(k)}(\mathbf{T})] - \tilde{L}^{(k)*}$.

Then there exist constants $A^{(k)}, S^{(k)} > 0$ such that $\mathbb{E}_{\mu_s^{(k)}}[\tilde{L}^{(k)}(\mathbf{T})] - \tilde{L}^{(k)*} \leq A^{(k)}s$ for all $s \in (0, S^{(k)})$. If we further choose the learning rate s such that $s \leq \min\{\frac{\epsilon}{2A^{(k)}}, S^{(k)}\}$, and the time t satisfies, $t \geq \frac{1}{\lambda_s^{(k)}} \log\left(\frac{2D^{(k)}(s, p_0)}{\epsilon}\right)$, where $D^{(k)}(s, p_0) = C^{(k)}(s) \cdot \|p_0 - \mu_s^{(k)}\|_{L^2((\mu_s^{(k)})^{-1})}$ and $C^{(k)}(s)$ is a positive constant, then,

$$\mathbb{E}[\tilde{L}^{(k)}(\mathbf{T}_t)] - \tilde{L}^{(k)*} \leq \epsilon. \quad (9)$$

The proofs of the above are given in Section A.1, Appendix C and Appendix E, respectively.

We note that since λ and ϵ can be set to be arbitrarily small positive numbers for the above convergence, it follows that such a mild regularizer would have negligible effect at small/finite weight values w.r.t unregularized loss and that the regularization only appreciably affects the shape of the loss at infinity.

Remark 3.4 (Necessity of Factor Regularization). We note that the factorized loss function $\mathcal{L}(\mathbf{T}) = \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{U}\mathbf{V}^\top)$, where $\mathbf{T} = (\mathbf{U}, \mathbf{V})$, exhibits a scaling invariance under the transformation, $g_{\mathbf{A}}(\mathbf{U}, \mathbf{V}) = (\mathbf{U}\mathbf{A}, \mathbf{V}\mathbf{A}^{-1})$, $\mathbf{A} \in \mathbb{R}^{r \times r}$, $\det(\mathbf{A}) \neq 0$, since,

$$\mathcal{L}(\mathbf{U}\mathbf{A}, \mathbf{V}\mathbf{A}^{-1}) = \mathcal{L}(\mathbf{U}, \mathbf{V}).$$

So the potential $\mathcal{L}(\mathbf{T})$ is constant along the non-compact orbits generated by the general linear group.

As a consequence, the Gibbs' measure $\mu_\gamma \propto e^{-\gamma\mathcal{L}(\mathbf{T})}$ is non-normalizable. Specifically, for any fixed rank- r matrix $\mathbf{W}_0 \neq 0$, $\mathcal{L}(\mathbf{T})$ is constant along the orbit $\mathcal{O} = \{(\mathbf{U}\mathbf{A}, \mathbf{V}\mathbf{A}^{-1}) \mid \det(\mathbf{A}) \neq 0\}$, which extends infinitely far from the origin. Hence, the partition function

$$Z = \int_{\mathbb{R}^D} e^{-\gamma\mathcal{L}(\mathbf{T})} d\mathbf{T}$$

diverges because it includes an integral of a non-zero constant density over an infinite-volume set. Consequently, the confining condition is violated, and the Poincaré Inequality cannot hold for the unregularized factorized loss.

4. An Empirical Study of Regularized Learning of Key and Query Matrices

Towards demonstrating an use of our regularized attention losses, we study the two-dimensional Darcy Flow PDE — which is popularly used as a benchmark in scientific-ML,

$$-\nabla \cdot (a(\mathbf{x}) \nabla u(\mathbf{x})) = f(\mathbf{x}), \quad \mathbf{x} \in (0, 1)^2, \quad (10)$$

where a is a spatially varying permeability (diffusion) coefficient, u is the unknown pressure field and f is a source function. The regression task is to learn a mapping $a \mapsto u$ from discretised input fields to discretised solution fields. Following the benchmark introduced by (Li et al., 2021), fields are discretised on a uniform 64×64 grid (bilinear-downsampled from the native 421×421 resolution). We use $N_{\text{train}} = 900$ samples for training and $N_{\text{test}} = 124$ held-out samples for evaluation. Both input and output fields are independently standardised using training-set mean and standard deviation.

Model Architecture

We use a patch-based single-head attention regressor defined as follows.

- (Tokenisation) Each 64×64 input field a is divided into non-overlapping 4×4 patches \mathbf{p}_i , yielding $t = (64/4)^2 = 256$ patches/tokens.
- (Embedding) Patches are projected to a d -dimensional representation via a two-stage convolutional encode: $\mathbf{x}_i = \text{ConvEncoder}(\mathbf{p}_i)$, and then a flattening layer. \mathbf{X} is then given by $[\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_t]^\top \in \mathbb{R}^{t \times d}$.
- (Positional encoding) Learnable 2-D positional bias \mathbf{P} for conv-token grid is added to the patch embeddings: $\hat{\mathbf{X}} = \mathbf{X} + \mathbf{P}$.
- (Single-head attention) Query, key, and value projections are

$$\begin{aligned} \mathbf{Q} &= \hat{\mathbf{X}} \mathbf{W}_Q, & \mathbf{K} &= \hat{\mathbf{X}} \mathbf{W}_K, & \mathbf{V} &= \hat{\mathbf{X}} \mathbf{W}_V, \\ \mathbf{W}_Q, \mathbf{W}_K &\in \mathbb{R}^{d \times r}, & \mathbf{W}_V &\in \mathbb{R}^{d \times d}. \end{aligned}$$

The attention map is computed with temperature fixed to 1 ($\beta = 1$), $\mathbf{A} = \text{RowSoftMax}_\beta\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right) \in \mathbb{R}^{t \times t}$.

The representation $\text{Attention}(\mathbf{X}) = \mathbf{A}\mathbf{V} \in \mathbb{R}^{t \times d}$ is projected back to patch space by a two-layer net and the resulting patches are rearranged into the predicted 64×64 output field \hat{u} . All dimensions are set to $r = d = 64$.

Two-Phase Training Protocol

Phase 1 (full pretraining): All parameters, embedding layers, projections \mathbf{W}_Q , \mathbf{W}_K , \mathbf{W}_V , and output MLP, are jointly optimised with plain mean-squared error (MSE) for 500 epochs using Adam (Kingma & Ba, 2014) with learning rate $\eta = 10^{-3}$ and batch size 32. This phase yields the *task-optimal* values for all weight matrices.

Phase 2 (ablation on ways of regularized training of \mathbf{W}_Q and \mathbf{W}_K): The goal is to isolate the effect of norm regularisation on the query and key matrices. Apart from the key and the query weight all other parameters are set to their Phase 1 values. Then \mathbf{W}_Q and \mathbf{W}_K are *re-initialised* identically across all three runs to their initial values used in Phase 1. Three objectives are then compared over 100 additional epochs,

1. **Unregularised:** (i.e., *none* in plots)

$$\mathcal{L}_0 = \hat{R}_*(\mathbf{T}),$$

2. **Log-amplified norm penalty:** (i.e., *log* in plots)

$$\begin{aligned} \mathcal{L}_{\log} &= \hat{R}_*(\mathbf{T}) + \frac{\lambda}{2} S \log(1 + S), \quad \lambda > 0, \\ S &= \|\mathbf{W}_Q\|_F^2 + \|\mathbf{W}_K\|_F^2. \end{aligned}$$

3. **Super-quadratic norm penalty:** (i.e., *power* in plots)

$$\mathcal{L}_{2+\epsilon} = \hat{R}_*(\mathbf{T}) + \frac{\lambda}{2} (\|\mathbf{W}_Q\|_F^{2+\epsilon} + \|\mathbf{W}_K\|_F^{2+\epsilon}), \quad \lambda, \epsilon > 0,$$

where $\mathbf{T} = (\mathbf{W}_Q, \mathbf{W}_K)$ and $\hat{R}_*(\mathbf{T})$ denotes the MSE loss with embedding layers, \mathbf{W}_V and output MLP frozen at their optimal values.

In all Phase 2 runs we use Adam with $\eta = 10^{-3}$, batch size 32 (as in Phase-1), and hyper-parameters $\epsilon = 10^{-6}$, $\lambda = 10^{-5}$ for *log* and $\lambda = 10^{-4}$ for *power*.

Metrics

We track the following quantities per epoch and per run.

- **Train/Test RMSE:** $\sqrt{\mathbb{E}[\|\hat{u} - u\|^2]}$ in normalised target space.
- **Test Relative L2 Error:** $\varepsilon_{\text{rel}} = \frac{\|\hat{u} - u\|_2}{\|u\|_2}$, the standard PDE-operator-learning benchmark metric.

- **Q/K Norm²:** $S = \|\mathbf{W}_Q\|_F^2 + \|\mathbf{W}_K\|_F^2$.

- **Generalisation gap:** $\Delta = \text{Test MSE} - \text{Train MSE}$.

4.1. Results analysis

In the very small regularization regime ($\epsilon = 10^{-6}$ and $\lambda = 10^{-4}$ or 10^{-5}), all methods, no regularization, “log” regularization, and “power” regularisation, exhibit nearly identical convergence behaviour in both normalised test/train RMSE and test relative L2 error. As shown in Figure 1, all configurations rapidly decay from their initial error and stabilise at approximately the same performance level (0.78), indicating that at this scale of regularization, there is no meaningful difference in predictive accuracy.

Despite the similarity in error metrics, substantial differences emerge in the internal dynamics of the model. In particular, the Frobenius norm ($\|\mathbf{W}_Q\|_F^2 + \|\mathbf{W}_K\|_F^2$) grows steadily in the absence of regularization, reaching significantly larger values over training. In contrast, both “log” and “power” regularization effectively constrain this growth, with power regularization enforcing the strongest suppression and log regularization yielding a slightly higher but still stable norm plateau. *Firstly*, these results demonstrate that even weak regularization of the kind studied here can meaningfully alter the scaling of the attention mechanism without impacting predictive performance.

In the generalization gap plot in Figure 1, all methods initially exhibit a rapid reduction in gap but the unregularized model shows a gradual increase over time, suggesting mild overfitting. Both log and power regularization mitigate this effect, maintaining a consistently lower and more stable gap throughout training, with power again providing the strongest control. *Secondly*, this result indicates that despite comparable test error improved generalization stability is induced by the regularizers being studied.

Overall, these findings suggest that the regularization studied in the presented theory primarily influences model stability. In low- λ regimes, log regularisation provides a robust default by stabilising the key and query weight norms with minimal intervention, while power regularisation is preferable when stronger suppression of attention weight growth is desired.

5. Conclusion

This work establishes that the SDE, that mimics the SGD, converges for both attention layers and depth-2 neural networks trained with LoRA, for arbitrary data and network sizes, even under arbitrarily low regularization, and we further provides non-asymptotic convergence rates.

A natural next question is whether the loss function on the attention layer with its key-query-value matrices using LoRA, satisfy the Villani conditions in the corresponding

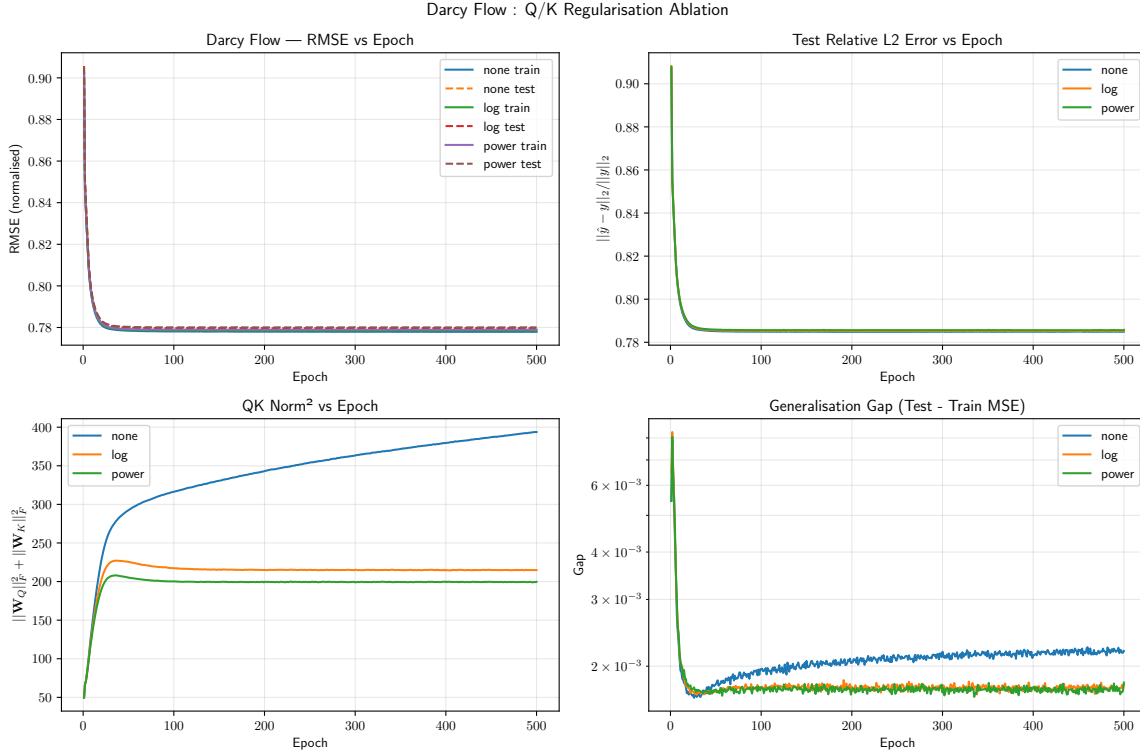


Figure 1. Per-epoch and per-run statistics for different loss functions

space of factor matrices. More generally, it remains open to prove convergence guarantees for a full transformer layer in which all three matrices per head (\mathbf{W}_Q , \mathbf{W}_K , \mathbf{W}_V) and the feedforward network weights are trained jointly, with one or more of these components potentially using LoRA.

An interesting direction for future work is to extend our guarantees to more advanced sequence models such as FlashAttention (Dao et al., 2022), Performers (Choromanski et al., 2021b), Mamba (Gu & Dao, 2023). In contrast to vanilla attention, which computes softmax attention, these methods modify the computation in different ways. FlashAttention (Dao et al., 2022) computes the same softmax attention but using a memory-efficient tiled algorithm, which retains the same mathematical structure as vanilla attention and is thus likely amenable to a similar analysis. Performers (Choromanski et al., 2021a) approximates the softmax kernel $\exp(\mathbf{Q}\mathbf{K}^\top)$ using positive random feature maps $\phi(\mathbf{Q})$ and $\phi(\mathbf{K})$, resulting in an approximate attention computation of the form $\text{Attn}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \approx \frac{\phi(\mathbf{Q})(\phi(\mathbf{K})^\top \mathbf{V})}{\phi(\mathbf{Q})(\phi(\mathbf{K})^\top \mathbf{1})}$, which reduces the computational complexity. Mamba (Gu & Dao, 2023) replaces attention entirely with an input-dependent state-space model (SSM), yielding a linear-time recurrence. Extending convergence guarantees to these alternative attention mechanisms remains an open problem.

Impact Statement

This paper presents work whose goal is to advance the theoretical underpinnings of the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Aghajanyan, A., Gupta, S., and Zettlemoyer, L. Intrinsic dimensionality explains the effectiveness of language model fine-tuning. In *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing (volume 1: long papers)*, pp. 7319–7328, 2021.
- Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate, 2016. URL <https://arxiv.org/abs/1409.0473>.
- Beltagy, I., Peters, M. E., and Cohan, A. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., and Tian, Q. Accurate medium-range global weather forecasting with 3d neural networks. *Nature*, 619:533–538, 2023.

- 440 Chen, S. and Li, Y. Provably learning a multi-head
441 attention layer. In *Proceedings of the 57th Annual*
442 *ACM Symposium on Theory of Computing, STOC '25*,
443 pp. 1744–1754, New York, NY, USA, 2025. Association
444 for Computing Machinery. ISBN 9798400715105.
445 doi: 10.1145/3717823.3718174. URL [https://doi.](https://doi.org/10.1145/3717823.3718174)
446 [org/10.1145/3717823.3718174](https://doi.org/10.1145/3717823.3718174).
- 447 Chen, S., Sheen, H., Wang, T., and Yang, Z. Unveiling
448 induction heads: Provable training dynamics and fea-
449 ture learning in transformers. In Globerson, A., Mackey,
450 L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., and
451 Zhang, C. (eds.), *Advances in Neural Information Pro-*
452 *cessing Systems*, volume 37, pp. 66479–66567. Curran
453 Associates, Inc., 2024. doi: 10.52202/079017-2127.
- 454 Child, R., Gray, S., Radford, A., and Sutskever, I. Gen-
455 erating long sequences with sparse transformers. *arXiv*
456 *preprint arXiv:1904.10509*, 2019.
- 457 Choromanski, K., Likhoshesterov, V., Dohan, D., Song, X.,
458 Gane, A., Sarlos, T., Hawkins, P., Davis, J., Mohiuddin,
459 A., Kaiser, L., Belanger, D., Colwell, L., and Weller, A.
460 Rethinking attention with performers. *ICLR*, 2021a.
- 461 Choromanski, K., Likhoshesterov, V., Dohan, D., Song, X.,
462 Gane, A., Sarlos, T., Hawkins, P., Davis, J., Mohiuddin,
463 A., Kaiser, L., et al. Rethinking attention with performers.
464 *International Conference on Learning Representations*,
465 2021b.
- 466 Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q. V., and
467 Salakhutdinov, R. Transformer-xl: Attentive language
468 models beyond a fixed-length context. *ACL*, 2019.
- 469 Dao, T., Fu, D. Y., Ermon, S., Rudra, A., and Ré, C. Flashat-
470 tention: Fast and memory-efficient exact attention with
471 io-awareness. *Advances in Neural Information Process-*
472 *ing Systems*, 2022.
- 473 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn,
474 D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer,
475 M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby,
476 N. An image is worth 16x16 words: Transformers for
477 image recognition at scale. In *International Conference*
478 *on Learning Representations*, 2021. URL [https://](https://openreview.net/forum?id=YicbFdNTTy)
479 openreview.net/forum?id=YicbFdNTTy.
- 480 Gao, C., Cao, Y., Li, Z., He, Y., Wang, M., Liu, H., Klu-
481 sowski, J., and Fan, J. Global convergence in training
482 large-scale transformers. *Advances in Neural Information*
483 *Processing Systems*, 37:29213–29284, 2024.
- 484 Gopalani, P. and Mukherjee, A. Global convergence of
485 sgd on two layer neural nets. *Information and Inference:*
486 *A Journal of the IMA*, 14(1):iaae035, 01 2025. ISSN
487 2049-8772. doi: 10.1093/imaia/iaae035. URL [https:](https://doi.org/10.1093/imaia/iaae035)
488 [//doi.org/10.1093/imaia/iaae035](https://doi.org/10.1093/imaia/iaae035).
- 489 Gopalani, P., Jha, S., and Mukherjee, A. Global convergence
490 of SGD for logistic loss on two layer neural nets. *Trans-*
491 *actions on Machine Learning Research*, 2024. ISSN 2835-
492 8856. URL [https://openreview.net/forum?](https://openreview.net/forum?id=9TqAUYB6tC)
493 [id=9TqAUYB6tC](https://openreview.net/forum?id=9TqAUYB6tC).
- 494 Gu, A. and Dao, T. Mamba: Linear-time sequence
modeling with selective state spaces. *arXiv preprint*
arXiv:2312.00752, 2023.
- Herde, M., Raonić, B., Rohner, T., Käppli, R., Molinaro,
R., de Bézenac, E., and Mishra, S. Poseidon: Efficient
foundation models for pdes. In Globerson, A., Mackey,
L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., and
Zhang, C. (eds.), *Advances in Neural Information Pro-*
cessing Systems, volume 37, pp. 72525–72624. Curran
Associates, Inc., 2024. doi: 10.52202/079017-2311.
- Hron, J., Bahri, Y., Sohl-Dickstein, J., and Novak, R. In-
finite attention: NNGP and NTK for deep attention
networks. In III, H. D. and Singh, A. (eds.), *Pro-*
ceedings of the 37th International Conference on Ma-
chine Learning, volume 119 of *Proceedings of Machine*
Learning Research, pp. 4376–4386. PMLR, 13–18 Jul
2020. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v119/hron20a.html)
[v119/hron20a.html](https://proceedings.mlr.press/v119/hron20a.html).
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang,
S., Wang, L., Chen, W., et al. Lora: Low-rank adaptation
of large language models. *Iclr*, 1(2):3, 2022.
- Jang, U., Lee, J. D., and Ryu, E. K. Lora training in the ntk
regime has no spurious local minima. In *International*
Conference on Machine Learning (ICML), 2024.
- Kim, J., Kim, J., and Ryu, E. K. Lora training provably
converges to a low-rank global minimum or it fails loudly.
arXiv preprint arXiv:2502.09376, 2025.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic
optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kumar, D., Jha, S., and Mukherjee, A. Langevin monte-
carlo provably learns depth two neural nets at any size
and data, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2503.10428)
[2503.10428](https://arxiv.org/abs/2503.10428).
- Lam, R., Sanchez-Gonzalez, A., Willson, M., Wyrnsberger,
P., Fortunato, M., Alet, F., Ravuri, S., Ewalds, T., Eaton-
Rosen, Z., Hu, W., et al. Learning skillful medium-range
global weather forecasting. *Science*, 382(6677):1416–
1421, 2023.
- Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P.,
and Soricut, R. Albert: A lite bert for self-supervised
learning of language representations. *ICLR*, 2020.

- 495 Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhat-
 496 tacharya, K., Stuart, A., and Anandkumar, A. Fourier neu-
 497 ral operator for parametric partial differential equations.
 498 *International Conference on Learning Representations*
 499 (*ICLR*), 2021.
- 500 Lin, L., Bai, Y., and Mei, S. Transformers as decision
 501 makers: Provable in-context reinforcement learning via
 502 supervised pretraining. *arXiv preprint arXiv:2310.08566*,
 503 2023.
- 504 Luong, T., Pham, H., and Manning, C. D. Effective ap-
 505 proaches to attention-based neural machine translation.
 506 In Márquez, L., Callison-Burch, C., and Su, J. (eds.),
 507 *Proceedings of the 2015 Conference on Empirical Meth-*
 508 *ods in Natural Language Processing*, pp. 1412–1421,
 509 Lisbon, Portugal, September 2015. Association for Com-
 510 putational Linguistics. doi: 10.18653/v1/D15-1166. URL
 511 <https://aclanthology.org/D15-1166/>.
 512
- 513 Molinaro, R., Lanthaler, S., Raonič, B., Rohner, T., Arme-
 514 gioiu, V., Simonis, S., Grund, D., Ramic, Y., Wan, Z. Y.,
 515 Sha, F., Mishra, S., and Zepeda-Núñez, L. Generative
 516 ai for fast and accurate statistical computation of fluids.
 517 *arXiv preprint arXiv:2409.18359*, 2024.
- 518 Ormaniec, W., Dangel, F., and Singh, S. P. What does it
 519 mean to be a transformer? insights from a theoretical
 520 hessian analysis. *arXiv preprint arXiv:2410.10986*, 2024.
- 521 Pathak, J., Subramanian, S., Harrington, P., Raja, S.,
 522 Chattopadhyay, A., Mardani, M., Kurth, T., Hall, D.,
 523 Li, Z., Azizzadenesheli, K., et al. Fourcastnet: A
 524 global data-driven high-resolution weather model us-
 525 ing adaptive fourier neural operators. *arXiv preprint*
 526 *arXiv:2202.11214*, 2022.
- 527 Radford, A., Narasimhan, K., Salimans, T., and
 528 Sutskever, I. Improving language understanding
 529 by generative pre-training. 2018. URL [https://cdn.openai.com/research-covers/
 530 language-unsupervised/language_
 531 understanding_paper.pdf](https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf).
- 532 Shi, B., Su, W., and Jordan, M. I. On learning rates and
 533 schrödinger operators. *Journal of Machine Learning*
 534 *Research*, 24(379):1–53, 2023. URL [http://jmlr.
 535 org/papers/v24/20-364.html](http://jmlr.org/papers/v24/20-364.html).
- 536 Shuttleworth, R. S., Andreas, J., Torralba, A., and Sharma,
 537 P. LoRA vs full fine-tuning: An illusion of equiva-
 538 lence. In *The Thirty-ninth Annual Conference on Neural*
 539 *Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=xp7B8rkh7L>.
- 540 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,
 541 L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I.
 542 Attention is all you need. In Guyon, I., Luxburg, U. V.,
 543 Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S.,
 544 and Garnett, R. (eds.), *Advances in Neural Information*
 545 *Processing Systems*, volume 30. Curran Associates, Inc.,
 546 2017. URL [https://proceedings.neurips.
 547 cc/paper_files/paper/2017/file/
 548 3f5ee243547dee91fbd053c1c4a845aa-Paper.
 549 pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- Villani, C. *Hypocoercivity*, volume 202. American Mathe-
 matical Society, 2009.
- Wang, S., Li, B. Z., Khabsa, M., Fang, H., and Ma, H.
 Linformer: Self-attention with linear complexity. *arXiv*
preprint arXiv:2006.04768, 2020.
- Yun, C., Bhojanapalli, S., Rawat, A. S., Reddi, S., and
 Kumar, S. Are transformers universal approximators of
 sequence-to-sequence functions? In *International Confer-*
ence on Learning Representations, 2020. URL <https://openreview.net/forum?id=ByxRM0Ntvr>.

A. Proof of Villani Conditions for Regression on Attention

Towards stating the proofs we note the following notations,

Definition A.1 (Defining $\mathbf{Y}_i, \hat{\mathbf{Y}}_i, \mathbf{S}_i$ and \mathbf{E}_i). For a choice of training data as, $\{(\mathbf{X}_i, \mathbf{Y}_i) \in \mathbb{R}^{t \times d} \times \mathbb{R}^{t \times d} \mid i = 1, \dots, n\}$ and $\hat{R}_A(\mathbf{T}) = \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{T})$, the loss is $\ell_i(\mathbf{T}) = \frac{1}{2} \|\mathbf{Y}_i - \hat{\mathbf{Y}}_i(\mathbf{T})\|_F^2$. The output is defined as $\hat{\mathbf{Y}}_i(\mathbf{T}) := \mathbf{S}_i(\mathbf{T}) \mathbf{X}_i \mathbf{W}_v$, where $\mathbf{S}_i(\mathbf{T}) := \text{RowSoftMax}_\beta(\mathbf{M}_i) \in \mathbb{R}^{t \times t}$ and $\mathbf{M}_i = \frac{1}{\sqrt{d}} \mathbf{X}_i \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}_i^\top \in \mathbb{R}^{t \times t}$. We also define the error as, $\mathbf{E}_i := \hat{\mathbf{Y}}_i(\mathbf{T}) - \mathbf{Y}_i \in \mathbb{R}^{t \times d}$

Lemma A.2. *The norm of RowSoftMax_β is bounded by $\|\mathbf{S}_i\|_F \leq \sqrt{t}$. The norm of the Jacobian and Hessian of RowSoftMax_β from Definition 2.3 are bounded by $\beta B_{s'}$ and $\beta^2 B_{s''}$, respectively. That is*

$$\|\mathrm{d}\mathbf{S}\|_F \leq B_{s'} \beta \|\mathrm{d}\mathbf{M}\|_F, \quad (11)$$

and

$$\|\mathrm{d}^2 \mathbf{S}_i\|_F \leq \beta^2 B_{s''} \|\mathrm{d}\mathbf{M}_i\|_F^2 \quad (12)$$

where $B_{s'} = 2$ and $B_{s''} = 6t^2$ are finite constants, and β is the constant temperature parameter introduced from Definition 2.3.

Lemma A.3. *The bound of gradient and laplacian of $\hat{R}_A(\mathbf{T})$ are given by*

$$\|\nabla_{\mathbf{T}} \hat{R}_A(\mathbf{T})\| = \sqrt{\|\nabla_{\mathbf{W}_Q} \hat{R}_A(\mathbf{T})\|_F^2 + \|\nabla_{\mathbf{W}_K} \hat{R}_A(\mathbf{T})\|_F^2} \leq \left[(\sqrt{t} B_x B_w + B_y) B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \right] \|\mathbf{T}\| \quad (13)$$

and

$$|\Delta_{\mathbf{T}} \hat{R}_A(\mathbf{T})| = \Delta_{\mathbf{W}_Q} \hat{R}_A(\mathbf{T}) + \Delta_{\mathbf{W}_K} \hat{R}_A(\mathbf{T}) \leq \frac{B_x^4}{d} \left(((\beta B_{s'}) B_x B_w)^2 + (\sqrt{t} B_x B_w + B_y) B_x B_w (\beta^2 B_{s''}) \right) \|\mathbf{T}\|^2. \quad (14)$$

The above lemmas are proved in Section B.

A.1. Proof of Theorem 3.1 for Loss in Definition 2.4

Proof. We note that,

$$V_{ATT}(\mathbf{T}) = \hat{R}_A(\mathbf{T}) + R(\mathbf{T}) = \hat{R}_A(\mathbf{T}) + \frac{\lambda}{2} \|\mathbf{T}\|^2 \log(1 + \|\mathbf{T}\|^2), \quad (15)$$

where the regularization is defined as $R(\mathbf{T}) = \frac{\lambda}{2} \|\mathbf{T}\|^2 \log(1 + \|\mathbf{T}\|^2)$.

We start with the gradient of the regularization term,

$$\nabla R(\mathbf{T}) = \lambda \mathbf{T} \log(1 + \|\mathbf{T}\|^2) + \lambda \mathbf{T} \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2}, \quad (16)$$

so

$$\|\nabla R(\mathbf{T})\| = \lambda \|\mathbf{T}\| \left(\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right) \quad (17)$$

Using the expression of $\|\nabla R(\mathbf{T})\|$ from above and the upperbound on $\|\nabla_{\mathbf{T}} \hat{R}_A(\mathbf{T})\|$ from Lemma A.3, we have,

$$\begin{aligned} \|\nabla V_{ATT}(\mathbf{T})\|^2 &= \|\nabla \hat{R}_A(\mathbf{T}) + \nabla R(\mathbf{T})\|^2 \geq \|\nabla R(\mathbf{T})\|^2 - 2\|\nabla R(\mathbf{T})\| \sup \|\nabla_{\mathbf{T}} \hat{R}_A(\mathbf{T})\| \\ &\geq \lambda^2 \|\mathbf{T}\|^2 \left(\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right)^2 \\ &\quad - 2\lambda \beta \|\mathbf{T}\| \left(\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right) \cdot \left[(\sqrt{t} B_x B_w + B_y) B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \right] \|\mathbf{T}\| \end{aligned} \quad (18)$$

²For a matrix-valued function $\mathbf{S}(\mathbf{M})$, the first-order differential $\mathrm{d}\mathbf{S}$ represents the linear principal part of the change in \mathbf{S} given an infinitesimal perturbation $\mathrm{d}\mathbf{M}$, defined via the Taylor expansion: $\mathbf{S}(\mathbf{M} + \mathrm{d}\mathbf{M}) = \mathbf{S}(\mathbf{M}) + \mathrm{d}\mathbf{S} + \mathcal{O}(\|\mathrm{d}\mathbf{M}\|^2)$. Similarly, the second-order differential $\mathrm{d}^2 \mathbf{S}$ represents the quadratic variation, such that $\mathbf{S}(\mathbf{M} + \mathrm{d}\mathbf{M}) = \mathbf{S}(\mathbf{M}) + \mathrm{d}\mathbf{S} + \frac{1}{2} \mathrm{d}^2 \mathbf{S} + \mathcal{O}(\|\mathrm{d}\mathbf{M}\|^3)$, encapsulating the action of the Hessian tensor.

Next we note that,

$$\begin{aligned}
 \Delta R(\mathbf{T}) &= \sum_{k=1}^D \frac{\partial}{\partial T_k} \left(\lambda T_k \left[\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right] \right) \\
 &= D\lambda \left[\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right] + \sum_{k=1}^D \lambda T_k \left[\frac{2T_k}{1 + \|\mathbf{T}\|^2} + \frac{2T_k}{(1 + \|\mathbf{T}\|^2)^2} \right] \\
 &= D\lambda \left[\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right] + \frac{2\lambda\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} + \frac{2\lambda\|\mathbf{T}\|^2}{(1 + \|\mathbf{T}\|^2)^2}.
 \end{aligned} \tag{19}$$

Also recall the upper bound of $|\Delta_T \hat{R}_A(\mathbf{T})|$

$$\frac{1}{s} \|\nabla V_{ATT}\|^2 - \Delta V_{ATT} \geq \frac{1}{s} \|\nabla V_{ATT}\|^2 - |\Delta_T \hat{R}_A| - \Delta R(T)$$

given Lemma A.3, substitute equation 18 and 19 into Villani condition, we have,

$$\begin{aligned}
 \lim_{\|\mathbf{T}\| \rightarrow \infty} \left(\frac{1}{s} \|\nabla V_{ATT}\|^2 - \Delta V_{ATT} \right) &\geq \lim_{\|\mathbf{T}\| \rightarrow \infty} \frac{1}{s} \|\nabla V_{ATT}\|^2 - |\Delta_T \hat{R}_A| - \Delta R(T) \\
 &= \lim_{\|\mathbf{T}\| \rightarrow \infty} \|\mathbf{T}\|^2 \left[\underbrace{\frac{\lambda^2}{s} \left(\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right)^2}_{\text{from equation 18 } \rightarrow +\infty \text{ at } O(\log \|\mathbf{T}\|^2)} \right. \\
 &\quad \left. - \underbrace{\frac{2\lambda\beta}{s} \left[(\sqrt{t}B_x B_w + B_y) B_w B_x B_{s'} \frac{B_x^2}{\sqrt{d}} \right]}_{\text{from equation 18 } \rightarrow +\infty \text{ at } O(\log \|\mathbf{T}\|)} \left(\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right) \right. \\
 &\quad \left. - \underbrace{\beta^2 \frac{B_x^4}{d} \left((B_{s'} B_x B_w)^2 + (\sqrt{t}B_x B_w + B_y) B_x B_w B_{s''} \right)}_{\text{from equation 14 constant}} - \underbrace{\frac{D\lambda \left[\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right] + \frac{2\lambda\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} + \frac{2\lambda\|\mathbf{T}\|^2}{(1 + \|\mathbf{T}\|^2)^2}}_{\text{from equation 19 } \rightarrow 0} \frac{\|\mathbf{T}\|^2}{\|\mathbf{T}\|^2} \right].
 \end{aligned} \tag{20}$$

So Villani condition is satisfied for all β and r (which is contained in D), this is because the leading term from the gradient of regularization dominates, which is independent on β and r . \square

A.2. Proof of Theorem 3.1 for Loss in Definition 2.5

Proof. We note that,

$$V_{ATT,\varepsilon}(\mathbf{T}) = \hat{R}_{A,\varepsilon}(\mathbf{T}) + R_\varepsilon(\mathbf{T}) = \hat{R}_A(\mathbf{T}) + \left[\frac{\lambda}{2} (\|\mathbf{W}_Q\|_F^{2+\varepsilon} + \|\mathbf{W}_K\|_F^{2+\varepsilon}) \right], \tag{21}$$

where the regularization term is defined as $R_\varepsilon(\mathbf{T}) = \left[\frac{\lambda}{2} (\|\mathbf{W}_Q\|_F^{2+\varepsilon} + \|\mathbf{W}_K\|_F^{2+\varepsilon}) \right]$. Since we have,

$$\begin{aligned}
 \nabla R_\varepsilon(\mathbf{T}) &= \nabla \left[\frac{\lambda}{2} (\|\mathbf{W}_Q\|_F^{2+\varepsilon} + \|\mathbf{W}_K\|_F^{2+\varepsilon}) \right] = \left(\nabla_{\mathbf{W}_Q} \left[\frac{\lambda}{2} (\|\mathbf{W}_Q\|_F^2)^{\frac{2+\varepsilon}{2}} \right], \nabla_{\mathbf{W}_K} \left[\frac{\lambda}{2} (\|\mathbf{W}_K\|_F^2)^{\frac{2+\varepsilon}{2}} \right] \right) \\
 &= \frac{\lambda}{2} \left(\frac{2+\varepsilon}{2} (\|\mathbf{W}_Q\|_F^2)^{\frac{\varepsilon}{2}} \cdot \nabla_{\mathbf{W}_Q} (\|\mathbf{W}_Q\|_F^2), \frac{2+\varepsilon}{2} (\|\mathbf{W}_K\|_F^2)^{\frac{\varepsilon}{2}} \cdot \nabla_{\mathbf{W}_K} (\|\mathbf{W}_K\|_F^2) \right) \\
 &= \frac{\lambda}{2} \left(\frac{2+\varepsilon}{2} \|\mathbf{W}_Q\|_F^\varepsilon \cdot (2\mathbf{W}_Q), \frac{2+\varepsilon}{2} \|\mathbf{W}_K\|_F^\varepsilon \cdot (2\mathbf{W}_K) \right) \\
 &= \frac{\lambda}{2} \left((2+\varepsilon) \|\mathbf{W}_Q\|_F^\varepsilon \mathbf{W}_Q, (2+\varepsilon) \|\mathbf{W}_K\|_F^\varepsilon \mathbf{W}_K \right)
 \end{aligned} \tag{22}$$

so

$$\|\nabla R_\varepsilon(\mathbf{T})\|^2 = \left(\frac{\lambda}{2} \right)^2 (2+\varepsilon)^2 (\|\mathbf{W}_Q\|_F^{2+2\varepsilon} + \|\mathbf{W}_K\|_F^{2+2\varepsilon}) \geq \left(\frac{\lambda}{2} \right)^2 (2+\varepsilon)^2 2^{-\varepsilon} \|\mathbf{T}\|^{2+2\varepsilon}, \tag{23}$$

where we have used, $\frac{(\|\mathbf{W}_Q\|^2)^{1+\epsilon} + (\|\mathbf{W}_K\|^2)^{1+\epsilon}}{2} \geq \left(\frac{\|\mathbf{W}_Q\|^2 + \|\mathbf{W}_K\|^2}{2}\right)^{1+\epsilon} = 2^{-1-\epsilon} \|\mathbf{T}\|^{2+2\epsilon}$ by Jensen Inequality.

Since $\frac{\partial}{\partial \mathbf{W}_{Q,ij}} \|\mathbf{W}_Q\|_F^{2+\epsilon} = (2+\epsilon) \|\mathbf{W}_Q\|_F^\epsilon \mathbf{W}_{Q,ij}$ from above, we have

$$\frac{\partial^2}{\partial \mathbf{W}_{Q,ij}^2} \|\mathbf{W}_Q\|_F^{2+\epsilon} = (2+\epsilon)\epsilon \|\mathbf{W}_Q\|_F^{\epsilon-2} W_{Q,ij}^2 + (2+\epsilon) \|\mathbf{W}_Q\|_F^\epsilon. \quad (24)$$

By summing over all $d \times r$ elements of \mathbf{W}_Q :

$$\begin{aligned} \Delta_{\mathbf{W}_Q} \|\mathbf{W}_Q\|_F^{2+\epsilon} &= (2+\epsilon)\epsilon \|\mathbf{W}_Q\|_F^{\epsilon-2} \cdot \left(\sum_{i=1}^d \sum_{j=1}^r W_{Q,ij}^2 \right) + \sum_{i=1}^d \sum_{j=1}^r [(2+\epsilon) \|\mathbf{W}_Q\|_F^\epsilon] \\ &= (2+\epsilon) \|\mathbf{W}_Q\|_F^\epsilon \cdot \epsilon + (2+\epsilon) \|\mathbf{W}_Q\|_F^\epsilon \cdot dr = (2+\epsilon)(\epsilon + dr) \|\mathbf{W}_Q\|_F^\epsilon. \end{aligned} \quad (25)$$

By symmetry, the Laplacian for $\mathbf{W}_K \in \mathbb{R}^{d \times r}$ follows the exact same form: $\Delta_{\mathbf{W}_K} \|\mathbf{W}_K\|_F^{2+\epsilon} = (2+\epsilon)(\epsilon + dr) \|\mathbf{W}_K\|_F^\epsilon$. We have

$$\Delta R_\epsilon(\mathbf{T}) = \frac{\lambda}{2} (2+\epsilon) [(\epsilon + dr) \|\mathbf{W}_Q\|_F^\epsilon + (\epsilon + dr) \|\mathbf{W}_K\|_F^\epsilon] \leq \frac{\lambda}{2} (2+\epsilon)(2\epsilon + D) \|\mathbf{T}\|^\epsilon, \quad (26)$$

where $D = 2dr$.

For the potential, together with the upper bound of $\|\nabla \hat{R}_A\|$ from Lemma A.3 we have,

$$\begin{aligned} \|\nabla V_{ATT,\epsilon}\|^2 &= \|\nabla \hat{R}_A + \nabla R_\epsilon\|^2 \geq \|\nabla R_\epsilon\|^2 - 2\|\nabla R_\epsilon\| \|\nabla \hat{R}_A\| \\ &\geq \left(\frac{\lambda}{2}\right)^2 (2+\epsilon)^2 2^{-\epsilon} \|\mathbf{T}\|^{2+2\epsilon} - 2\left(\frac{\lambda}{2}\right) (2+\epsilon) \|\mathbf{T}\|^{1+\epsilon} \left(\left[(\sqrt{t} B_x B_w + B_y) B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \right] \|\mathbf{T}\| \right) \\ &= \left(\frac{\lambda}{2}\right)^2 (2+\epsilon)^2 2^{-\epsilon} \|\mathbf{T}\|^{2+2\epsilon} - \lambda(2+\epsilon) \left[(\sqrt{t} B_x B_w + B_y) B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \right] \|\mathbf{T}\|^{2+\epsilon}. \end{aligned} \quad (27)$$

The last inequality we have used $\|\mathbf{W}_Q\|_F^{2+2\epsilon} + \|\mathbf{W}_K\|_F^{2+2\epsilon} \leq (\|\mathbf{W}_Q\|_F^2 + \|\mathbf{W}_K\|_F^2)^{1+\epsilon} = \|\mathbf{T}\|^{2+2\epsilon}$.

Substituting this, along with the upper bound of $|\Delta_T \hat{R}_A|$ by Lemma A.3 and upper bound of $\Delta R_\epsilon(\mathbf{T})$ from equation 26, into the Villani limit expression for any given $s > 0$,

$$\begin{aligned} \lim_{\|\mathbf{T}\| \rightarrow \infty} \left(\frac{1}{s} \|\nabla V_{ATT,\epsilon}\|^2 - \Delta V_{ATT,\epsilon} \right) &\geq \lim_{\|\mathbf{T}\| \rightarrow \infty} \frac{1}{s} \|\nabla V_{ATT,\epsilon}\|^2 - |\Delta_T \hat{R}_A| - \Delta R_\epsilon(\mathbf{T}) \\ &\geq \lim_{\|\mathbf{T}\| \rightarrow \infty} \left[\frac{1}{s} \left(\left(\frac{\lambda}{2}\right)^2 (2+\epsilon)^2 2^{-\epsilon} \|\mathbf{T}\|^{2+2\epsilon} - \lambda(2+\epsilon) \left[(\sqrt{t} B_x B_w + B_y) B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \right] \|\mathbf{T}\|^{2+\epsilon} \right) \right. \\ &\quad \left. - \left((\beta B_{s'} B_x B_w)^2 + (\sqrt{t} B_x B_w + B_y) B_x B_w \beta^2 B_{s''} \right) \frac{B_x^4}{d} \|\mathbf{T}\|^2 - \frac{\lambda}{2} (2+\epsilon)(2\epsilon + D) \|\mathbf{T}\|^\epsilon \right] \\ &= \lim_{\|\mathbf{T}\| \rightarrow \infty} \|\mathbf{T}\|^2 \left[\underbrace{\frac{1}{s} \left(\frac{\lambda}{2}\right)^2 (2+\epsilon)^2 2^{-\epsilon} \|\mathbf{T}\|^{2\epsilon}}_{\text{from equation 27} \rightarrow +\infty \text{ at } \|\mathbf{T}\|^{2\epsilon}} - \underbrace{\frac{\lambda(2+\epsilon) \left[(\sqrt{t} B_x B_w + B_y) B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \right]}{s}}_{\text{from equation 27} \rightarrow +\infty \text{ at } \|\mathbf{T}\|^\epsilon} \|\mathbf{T}\|^\epsilon \right. \\ &\quad \left. - \underbrace{\left((\beta B_{s'} B_x B_w)^2 + (\sqrt{t} B_x B_w + B_y) B_x B_w \beta^2 B_{s''} \right) \frac{B_x^4}{d}}_{\text{from equation 14} \text{ Constant}} - \underbrace{\frac{\lambda}{2} (2+\epsilon)(2\epsilon + D)}_{\text{from equation 26} \rightarrow 0} \|\mathbf{T}\|^{\epsilon-2} \right]. \end{aligned} \quad (28)$$

The leading order term $O(\|\mathbf{T}\|^{2\epsilon})$ dominates. Since $\epsilon > 0$, regardless of the choices of dimensions d, t , key-query inner dimension r , and Softmax temperature β . Consequently, the limit diverges to $+\infty$, proving that the ϵ -regularized attention loss unconditionally satisfies the Villani condition. \square

B. Proofs of Intermediate Lemmas for Theorem 3.1

Proof of Lemma A.2. By the definition of the row-wise softmax function, each element $S_{i,j,k}$ (the element in the j -th row and k -th column of S_i) represents a valid probability. Therefore, all elements are positive, $S_{i,j,k} \in (0, 1)$, and the sum of the elements across any given row j is exactly 1: $\sum_{k=1}^t S_{i,j,k} = 1$. The squared Frobenius norm of the matrix S_i is defined as the sum of its squared entries:

$$\|S_i\|_F^2 = \sum_{j=1}^t \sum_{k=1}^t S_{i,j,k}^2 \leq \left(\sum_{k=1}^t S_{i,j,k} \right)^2 = 1^2 = 1,$$

so $\|S_i\|_F^2 \leq \sum_{j=1}^t 1 = t$, which means the bound of $\mathbf{S} = \text{RowSoftMax}_\beta(\mathbf{M})$ is given by $\|S_i\|_F \leq \sqrt{t}$.

To determine the exact bounds for the Jacobian and Hessian of $\mathbf{S} = \text{RowSoftMax}_\beta(\mathbf{M})$, we adapt the structural findings from (Ormaniec et al., 2024).

As established in Appendix C.1, Lemma C.1 of (Ormaniec et al., 2024), since the row-wise softmax applies independently to each row, the cross-row derivatives are strictly zero. This decoupling endows the full Jacobian matrix and the Hessian tensor with a block-diagonal structure, allowing us to analyze them row by row.

Jacobian Bound For any single row i , the output probability is defined element-wise as $S_{i,j} = \frac{\exp(\beta M_{i,j})}{\sum_l \exp(\beta M_{i,l})}$. To compute the derivative with respect to the input $M_{i,k}$ we invoke Lemma B.1, Equation 19 of (Ormaniec et al., 2024) to obtain the exact local Jacobian matrix,

$$\frac{\partial S_{i,j}}{\partial M_{i,k}} = \frac{\partial}{\partial(\beta M_{i,k})} \left(\frac{\exp(\beta M_{i,j})}{\sum_l \exp(\beta M_{i,l})} \right) \cdot \frac{\partial(\beta M_{i,k})}{\partial M_{i,k}} = S_{i,j}(\delta_{j,k} - S_{i,k}) \cdot \beta \quad (29)$$

where $\delta_{j,k}$ is the Kronecker delta. Expressing this element-wise relationship in matrix form for the entire i -th row $\mathbf{S}_i \in \mathbb{R}^t$,

$$\mathcal{J}_i = \frac{\partial \mathbf{S}_i}{\partial \mathbf{M}_i} = \beta (\text{diag}(\mathbf{S}_i) - \mathbf{S}_i^\top \mathbf{S}_i) \quad (30)$$

where $\mathcal{J}_i \in \mathbb{R}^{t \times t}$ is used to denote the Jacobian and \mathbf{S}_i is the i -th row of the attention probability matrix. Since $\mathbf{S}_{i,k} \in (0, 1)$ and $\sum_k \mathbf{S}_{i,k} = 1$, $\text{diag}(\mathbf{S}_i) - \mathbf{S}_i^\top \mathbf{S}_i$ represents the exact covariance matrix of a categorical distribution. We have the bound

$$\|\text{diag}(\mathbf{S}_i) - \mathbf{S}_i^\top \mathbf{S}_i\|_2 \leq \|\text{diag}(\mathbf{S}_i)\|_2 + \|\mathbf{S}_i^\top \mathbf{S}_i\|_2 \leq \max_k \mathbf{S}_{i,k} + \sum_{k=1}^t \mathbf{S}_{i,k}^2 \leq 2. \quad (31)$$

In the above step, the first inequality applies the triangle inequality. The subsequent equality is exact for the spectral norm ($\|\cdot\|_2$): the norm of the diagonal matrix equals its maximum entry, and the norm of the rank-1 positive semi-definite matrix $\mathbf{S}_i^\top \mathbf{S}_i$ equals its trace. The final strict inequality holds because \mathbf{S}_i is a probability vector which dictates that $\max_k \mathbf{S}_{i,k} < 1$ and $\sum_{k=1}^t \mathbf{S}_{i,k}^2 < 1$.

We denote $d\mathbf{M}$ and $d\mathbf{S}$ as the first-order matrix differentials, representing an arbitrary infinitesimal perturbation in the input pre-activation matrix and the corresponding induced perturbation in the output probability matrix, respectively. For a single row i , $d\mathbf{M}_i$ and $d\mathbf{S}_i$ represent their respective row vector differentials. Since the full Jacobian is block-diagonal, the differential mapping from $d\mathbf{M}$ to $d\mathbf{S}$ operates independently on each row. For any individual row i , the Euclidean norm (2-norm) of the differential vector satisfies $\|d\mathbf{S}_i\|_2 \leq 2\beta \|d\mathbf{M}_i\|_2$ based on the spectral norm bound derived above. By definition, the squared Frobenius norm of a matrix is the sum of the squared 2-norms of its row vectors. Summing over all t rows, we obtain:

$$\|d\mathbf{S}\|_F^2 = \sum_{i=1}^t \|d\mathbf{S}_i\|_2^2 \leq \sum_{i=1}^t (2\beta)^2 \|d\mathbf{M}_i\|_2^2 = 4\beta^2 \sum_{i=1}^t \|d\mathbf{M}_i\|_2^2 = 4\beta^2 \|d\mathbf{M}\|_F^2 \quad (32)$$

Taking the square root of both sides directly yields the global bound for the differential:

$$\|d\mathbf{S}\|_F \leq B_{s'} \beta \|d\mathbf{M}\|_F \quad (33)$$

Therefore, the Jacobian norm is bounded by $\beta B_{s'}$, where the constant is explicitly evaluated as $B_{s'} = 2$.

Hessian Bound Differentiating the local Jacobian element $\frac{\partial \mathbf{S}_{i,j}}{\partial \mathbf{M}_{i,k}} = \beta(\mathbf{S}_{i,j}\delta_{j,k} - \mathbf{S}_{i,j}\mathbf{S}_{i,k})$ with respect to another input $\mathbf{M}_{i,l}$ yields an additional factor of β via the chain rule.

$$\begin{aligned} \frac{\partial^2 \mathbf{S}_{i,j}}{\partial \mathbf{M}_{i,k} \partial \mathbf{M}_{i,l}} &= \beta \left(\frac{\partial \mathbf{S}_{i,j}}{\partial \mathbf{M}_{i,l}} \delta_{j,k} - \frac{\partial \mathbf{S}_{i,j}}{\partial \mathbf{M}_{i,l}} \mathbf{S}_{i,k} - \mathbf{S}_{i,j} \frac{\partial \mathbf{S}_{i,k}}{\partial \mathbf{M}_{i,l}} \right) \\ &= \beta^2 ((\mathbf{S}_{i,j}\delta_{j,l} - \mathbf{S}_{i,j}\mathbf{S}_{i,l})\delta_{j,k} - (\mathbf{S}_{i,j}\delta_{j,l} - \mathbf{S}_{i,j}\mathbf{S}_{i,l})\mathbf{S}_{i,k} - \mathbf{S}_{i,j}(\mathbf{S}_{i,k}\delta_{k,l} - \mathbf{S}_{i,k}\mathbf{S}_{i,l})) \\ &= \beta^2 \mathbf{S}_{i,j} (\delta_{j,k}\delta_{j,l} - \mathbf{S}_{i,l}\delta_{j,k} - \delta_{j,l}\mathbf{S}_{i,k} + \mathbf{S}_{i,l}\mathbf{S}_{i,k} - \mathbf{S}_{i,k}\delta_{k,l} + \mathbf{S}_{i,k}\mathbf{S}_{i,l}) \\ &= \beta^2 \mathbf{S}_{i,j} (2\mathbf{S}_{i,k}\mathbf{S}_{i,l} + \delta_{j,k}\delta_{j,l} - \delta_{k,l}\mathbf{S}_{i,k} - \delta_{j,k}\mathbf{S}_{i,l} - \delta_{j,l}\mathbf{S}_{i,k}) \end{aligned} \quad (34)$$

This scalar element-wise formulation is the direct expansion of the matrix-level second derivative derived in Lemma C.1 of (Ormaniec et al., 2024). This expression exclusively comprises attention probability values $\mathbf{S}_{i,\cdot} \in (0, 1)$ and Kronecker deltas $\delta \in \{0, 1\}$. By the triangle inequality, the absolute value of each individual entry in this local 3D Hessian tensor $\mathcal{H}_i \in \mathbb{R}^{t \times t \times t}$, defined as the local 3D Hessian tensor $\mathcal{H}_i \in \mathbb{R}^{t \times t \times t}$ as the collection of all second-order partial derivatives of the i -th row of the output \mathbf{S} with respect to the i -th row of the input \mathbf{M} . For $i \in \{1, 2, \dots, t\}$, its element at index (j, k, l) is defined as:

$$(\mathcal{H}_i)_{j,k,l} := \frac{\partial^2 \mathbf{S}_{i,j}}{\partial \mathbf{M}_{i,k} \partial \mathbf{M}_{i,l}} \quad \text{for } j, k, l \in \{1, \dots, t\},$$

which is strictly bounded as,

$$\left| \frac{\partial^2 \mathbf{S}_{i,j}}{\partial \mathbf{M}_{i,k} \partial \mathbf{M}_{i,l}} \right| < \beta^2 \cdot 1 \cdot (2 \cdot 1 \cdot 1 + 1 + 1 + 1 + 1) = 6\beta^2 \quad (35)$$

For a single row i , the squared Frobenius norm of its local Hessian tensor \mathcal{H}_i is the sum of its t^3 squared entries. We bound this local tensor norm as,

$$\|\mathcal{H}_i\|_F^2 = \sum_{j,k,l=1}^t \left(\frac{\partial^2 \mathbf{S}_{i,j}}{\partial \mathbf{M}_{i,k} \partial \mathbf{M}_{i,l}} \right)^2 \leq \sum_{j,k,l=1}^t (6\beta^2)^2 = t^3 (36\beta^4) \implies \|\mathcal{H}_i\|_F \leq 6\beta^2 t^{1.5} \quad (36)$$

For mapping $\mathbf{S} = \text{RowSoftMax}_\beta(\mathbf{M})$, the complete global 6th-order Hessian tensor $H_{full} \in \mathbb{R}^{t \times t \times t \times t \times t \times t}$ is defined as the derivative of any output entry matrix $\mathbf{S}_{i,j}$ with respect to any two input entries $\mathbf{M}_{p,k}$ and $\mathbf{M}_{q,l}$:

$$(H_{full})_{i,j,p,k,q,l} := \frac{\partial^2 \mathbf{S}_{i,j}}{\partial \mathbf{M}_{p,k} \partial \mathbf{M}_{q,l}} \quad \text{for } i, j, p, k, q, l \in \{1, \dots, t\}. \quad (37)$$

As discussed in Appendix C.1, Lemma C.1 of (Ormaniec et al., 2024), the full Hessian tensor for the entire matrix-to-matrix mapping is block-diagonal, all cross-row second derivatives evaluate to zero. We have

$$\frac{\partial^2 \mathbf{S}_{i,j}}{\partial \mathbf{M}_{p,k} \partial \mathbf{M}_{q,l}} = \begin{cases} (\mathcal{H}_i)_{j,k,l} & \text{if } i = p = q, \\ 0 & \text{otherwise.} \end{cases} \quad (38)$$

This simplifies its total squared Frobenius norm being simply the sum of the squared norms of the t independent row blocks. Thus, the global Hessian norm $\|H_{full}\|_F$, satisfies,

$$\begin{aligned} \|H_{full}\|_F^2 &= \sum_{i=1}^t \sum_{j=1}^t \sum_{p=1}^t \sum_{k=1}^t \sum_{q=1}^t \sum_{l=1}^t \left(\frac{\partial^2 \mathbf{S}_{i,j}}{\partial \mathbf{M}_{p,k} \partial \mathbf{M}_{q,l}} \right)^2 \\ &= \sum_{i=1}^t \sum_{j,k,l=1}^t \left(\frac{\partial^2 \mathbf{S}_{i,j}}{\partial \mathbf{M}_{i,k} \partial \mathbf{M}_{i,l}} \right)^2 = \sum_{i=1}^t \|\mathcal{H}_i\|_F^2 \leq \sum_{i=1}^t (36\beta^4 t^3) = 36\beta^4 t^4. \end{aligned} \quad (39)$$

Taking the square root of both sides, the global Hessian norm is bounded by $6\beta^2 t^2$. Therefore, the bound can be denoted as $\beta^2 B_{s''}$, that is

$$\|d^2 \mathbf{S}_i\|_F \leq \beta^2 B_{s''} \|d\mathbf{M}_i\|_F^2, \quad (40)$$

where the constant is explicitly evaluated as $B_{s''} = 6t^2$. \square

Proof of Lemma A.3. We begin by bounding the norm of the gradient $\nabla_{\mathbf{T}} \hat{R}_A$, where $\mathbf{T} = (\mathbf{W}_Q, \mathbf{W}_K)$. Recall from Definitions 2.4 and 2.5 that $\hat{R}_A := \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{T})$. Taking gradients and applying the triangle inequality yields $\|\nabla_{\mathbf{T}} \hat{R}_A\| \leq \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{T}} \ell_i(\mathbf{T})\|$. Thus, it suffices to bound $\|\nabla_{\mathbf{T}} \ell_i(\mathbf{T})\|$, where $\|\nabla_{\mathbf{T}} \ell_i(\mathbf{T})\|^2 = \|\nabla_{\mathbf{W}_Q} \ell_i\|_F^2 + \|\nabla_{\mathbf{W}_K} \ell_i\|_F^2$.

From Definition A.1, recall that $M_i = \frac{1}{\sqrt{d}} \mathbf{X}_i \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}_i^\top$. In equation (54), we will show that the parameter gradients $\nabla_{\mathbf{W}_Q} \ell_i$ and $\nabla_{\mathbf{W}_K} \ell_i$ can be expressed in terms of the intermediate gradient $\nabla_{M_i} \ell_i$ as

$$\begin{aligned} \nabla_{\mathbf{W}_Q} \ell_i &= \frac{1}{\sqrt{d}} \mathbf{X}_i^\top (\nabla_{M_i} \ell_i) \mathbf{X}_i \mathbf{W}_K, \\ \nabla_{\mathbf{W}_K} \ell_i &= \frac{1}{\sqrt{d}} \mathbf{X}_i^\top (\nabla_{M_i} \ell_i) \mathbf{X}_i \mathbf{W}_Q. \end{aligned}$$

Therefore, it suffices to bound $\|\nabla_{M_i} \ell_i\|_F$. Next, we will express the intermediate gradient in terms of the upstream gradient $(\nabla_{S_i} \ell_i)$ in equation (48), $\nabla_{M_i} \ell_i = \left(\frac{\partial S_i}{\partial M_i} \right)^\top \nabla_{S_i} \ell_i$, where S_i is defined in Definition A.1. In equation (44) we will show that the upstream gradient admits the form $\nabla_{S_i} \ell_i = \mathbf{E}_i \mathbf{W}_v^\top \mathbf{X}_i^\top$, where $\mathbf{E}_i = \hat{\mathbf{Y}}_i - \mathbf{Y}_i$ is the error matrix.

Thus, we proceed in four steps, (Step 1) **Bounding the Error Matrix** (i.e. bounding $\|\mathbf{E}_i\|_F$), which in turn allows us to (Step 2) **Bounding the Upstream Gradient** (i.e. control $\|\nabla_{S_i} \ell_i\|_F$), then (Step 3) **Bounding the Intermediate Gradient** (i.e. $\|\nabla_{M_i} \ell_i\|_F$) and finally (Step 4) **Bounding the Parameter Gradients**.

• **Bounding the Error Matrix:** First using the triangle inequality and the Softmax output bound $\|S_i\|_F \leq \sqrt{t}$ to bound $\mathbf{E}_i = \hat{\mathbf{Y}}_i - \mathbf{Y}_i$ to obtain $\|\mathbf{E}_i\|_F \leq \sqrt{t} B_x B_w + B_y$ (see Equation 41).

Recalling $\mathbf{E}_i = \hat{\mathbf{Y}}_i(\mathbf{T}) - \mathbf{Y}_i$ we have,

$$\|\mathbf{E}_i\|_F \leq \|\hat{\mathbf{Y}}_i\|_F + \|\mathbf{Y}_i\|_F \leq \|S_i\|_F \|\mathbf{X}_i\|_F \|\mathbf{W}_v\|_F + B_y \leq \sqrt{t} B_x B_w + B_y, \quad (41)$$

where we have used the bound of $\|S_i\|_F$ proved in Lemma A.2. Recall the other constants used in the RHS come from definitions 2.3, 2.8 and 2.9.

• **Bounding the Upstream Gradient:** Based on the explicit expression for the loss gradient with respect to the output probability matrix, $\nabla_{S_i} \ell_i = \mathbf{E}_i \mathbf{W}_v^\top \mathbf{X}_i^\top$, we apply sub-multiplicativity to bound the upstream gradient: $\|\nabla_{S_i} \ell_i\|_F \leq (\sqrt{t} B_x B_w + B_y) B_w B_x$ (see Equation 45).

Recalling $\hat{\mathbf{Y}}_i(\mathbf{T}) = S_i(\mathbf{T}) \mathbf{X}_i \mathbf{W}_v$, for the loss $\ell_i(\mathbf{T}) = \frac{1}{2} \|\mathbf{Y}_i - \hat{\mathbf{Y}}_i(\mathbf{T})\|_F^2 = \frac{1}{2} \text{Tr}(\mathbf{E}_i^\top \mathbf{E}_i)$, take the differential

$$d\ell_i = \frac{1}{2} \text{Tr}((d\mathbf{E}_i)^\top \mathbf{E}_i + \mathbf{E}_i^\top d\mathbf{E}_i) = \text{Tr}(\mathbf{E}_i^\top d\mathbf{E}_i). \quad (42)$$

Since the target \mathbf{Y}_i is a constant, this is simply $d\mathbf{E}_i = d\hat{\mathbf{Y}}_i$. Recalling $\hat{\mathbf{Y}}_i(\mathbf{T}) = S_i(\mathbf{T}) \mathbf{X}_i \mathbf{W}_v$, \mathbf{X}_i and \mathbf{W}_v are treated as constants with respect to the pre-activations, we have $d\hat{\mathbf{Y}}_i = (dS_i) \mathbf{X}_i \mathbf{W}_v$. So the differential of the loss ℓ_i with respect to S_i is given by:

$$d\ell_i = \text{Tr}(\mathbf{E}_i^\top d\hat{\mathbf{Y}}_i) = \text{Tr}(\mathbf{E}_i^\top (dS_i) \mathbf{X}_i \mathbf{W}_v). \quad (43)$$

Using the cyclic property of the trace, we can rearrange this as $d\ell_i = \text{Tr}(\mathbf{X}_i \mathbf{W}_v \mathbf{E}_i^\top dS_i) = \text{Tr}((\mathbf{E}_i \mathbf{W}_v^\top \mathbf{X}_i^\top)^\top dS_i)$. By identifying this with the standard Frobenius inner product $d\ell_i = \text{Tr}((\nabla_{S_i} \ell_i)^\top dS_i)$, we extract the exact gradient,

$$\nabla_{S_i} \ell_i = \mathbf{E}_i \mathbf{W}_v^\top \mathbf{X}_i^\top. \quad (44)$$

So we have,

$$\|\nabla_{S_i} \ell_i\|_F = \|\mathbf{E}_i \mathbf{W}_v^\top \mathbf{X}_i^\top\|_F \leq \|\mathbf{E}_i\|_F \|\mathbf{W}_v\|_F \|\mathbf{X}_i\|_F \leq (\sqrt{t} B_x B_w + B_y) B_w B_x. \quad (45)$$

• **Bounding the Intermediate Gradient:** Utilizing the exact chain rule $\nabla_{M_i} \ell_i = \left(\frac{\partial S_i}{\partial M_i} \right)^\top \nabla_{S_i} \ell_i$ and substituting the Softmax Jacobian bound $\left\| \left(\frac{\partial S_i}{\partial M_i} \right) \right\|_F \leq \beta B_{s'}$, we obtain the intermediate gradient bound: $\|\nabla_{M_i} \ell_i\|_F \leq (\sqrt{t} B_x B_w + B_y) B_w B_x \beta B_{s'}$ (see Equation 49). Recall that the local Jacobian $\left(\frac{\partial S_i}{\partial M_i} \right) \in \mathbb{R}^{t \times t}$ and its bound are defined in Lemma A.2.

Since the loss differential can be equivalently expressed in terms of either M_i or S_i

$$d\ell_i = \langle \nabla_{M_i} \ell_i, dM_i \rangle = \langle \nabla_{S_i} \ell_i, dS_i \rangle, \quad (46)$$

express this in trace form, we have

$$\text{Tr}((\nabla_{M_i} \ell_i)^\top dM_i) = \text{Tr}((\nabla_{S_i} \ell_i)^\top dS_i) = \text{Tr}\left((\nabla_{S_i} \ell_i)^\top \left(\frac{\partial S_i}{\partial M_i}\right) dM_i\right) = \text{Tr}\left(\left(\left(\frac{\partial S_i}{\partial M_i}\right)^\top \nabla_{S_i} \ell_i\right)^\top dM_i\right), \quad (47)$$

we can extract

$$\nabla_{M_i} \ell_i = \left(\frac{\partial S_i}{\partial M_i}\right)^\top \nabla_{S_i} \ell_i. \quad (48)$$

Take Frobenius norm of both sides and apply the Cauchy-Schwarz inequality, we have

$$\|\nabla_{M_i} \ell_i\|_F \leq (\sqrt{t}B_x B_w + B_y) B_w B_x \beta B_{s'}. \quad (49)$$

• **Bounding the Parameter Gradients:** By extracting the gradient with respect to the query and key weight matrices, W_Q and W_K , we obtain the bound of $\|\nabla_{W_Q} \ell_i\|_F$ and $\|\nabla_{W_K} \ell_i\|_F$ using the bound of $\|\nabla_{M_i} \ell_i\|_F$, then we substitute the intermediate bound to achieve the final bound $\|\nabla_{\mathbf{T}} \ell_i\|^2 = \|\nabla_{W_Q} \ell_i\|_F^2 + \|\nabla_{W_K} \ell_i\|_F^2$, and hence the gradient squared bound (see Equation 58).

Next we fix W_k to analyze the partial derivative with respect to W_Q . To see that first consider the change in M_i due to variation in W_Q ,

$$d_{W_Q} M_i = \frac{1}{\sqrt{d}} X_i (dW_Q) W_K^\top X_i^\top, \quad (50)$$

Since,

$$d\ell_i = \langle \nabla_{M_i} \ell_i, d_{W_Q} M_i \rangle = \text{Tr}((\nabla_{M_i} \ell_i)^\top d_{W_Q} M_i) = \text{Tr}((\nabla_{M_i} \ell_i)^\top \frac{1}{\sqrt{d}} X_i (dW_Q) W_K^\top X_i^\top), \quad (51)$$

apply the cyclic property of trace, we have

$$d\ell_i = \frac{1}{\sqrt{d}} \text{Tr}(W_K^\top X_i^\top (\nabla_{M_i} \ell_i)^\top X_i (dW_Q)). \quad (52)$$

Compare this with the differential directly expressed in terms of the gradient with respect to W_Q as

$$d\ell_i = \langle \nabla_{W_Q} \ell_i, dW_Q \rangle = \text{Tr}((\nabla_{W_Q} \ell_i)^\top dW_Q), \quad (53)$$

we have $(\nabla_{W_Q} \ell_i)^\top = \frac{1}{\sqrt{d}} W_K^\top X_i^\top (\nabla_{M_i} \ell_i)^\top X_i$, or

$$\nabla_{W_Q} \ell_i = \frac{1}{\sqrt{d}} X_i^\top (\nabla_{M_i} \ell_i) X_i W_K. \quad (54)$$

By taking Frobenius norm of both sides and applying the Cauchy-Schwarz inequality, we have:

$$\|\nabla_{W_Q} \ell_i\|_F \leq \frac{1}{\sqrt{d}} \|X_i^\top\|_F \cdot \|\nabla_{M_i} \ell_i\|_F \cdot \|X_i\|_F \cdot \|W_K\|_F = \|\nabla_{M_i} \ell_i\|_F \frac{B_x^2}{\sqrt{d}} \|W_K\|_F. \quad (55)$$

Finally, substituting the upper bound $\|\nabla_{M_i} \ell_i\|_F$ derived from equation 49, we have:

$$\|\nabla_{W_Q} \ell_i\|_F \leq (\sqrt{t}B_x B_w + B_y) B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \|W_K\|_F. \quad (56)$$

By symmetry,

$$\|\nabla_{\mathbf{W}_K} \ell_i\|_F \leq (\sqrt{t}B_x B_w + B_y)B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \|\mathbf{W}_Q\|_F. \quad (57)$$

Combining equations 56 and 57, we have,

$$\|\nabla_{\mathbf{T}} \ell_i\|^2 = \|\nabla_{\mathbf{W}_Q} \ell_i\|_F^2 + \|\nabla_{\mathbf{W}_K} \ell_i\|_F^2 \leq \left[(\sqrt{t}B_x B_w + B_y)B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \right]^2 (\|\mathbf{W}_K\|_F^2 + \|\mathbf{W}_Q\|_F^2) \quad (58)$$

$$= \left[(\sqrt{t}B_x B_w + B_y)B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \right]^2 \|\mathbf{T}\|^2. \quad (59)$$

So for $\|\nabla_{\mathbf{T}} \hat{R}_A\|$, we have,

$$\|\nabla_{\mathbf{T}} \hat{R}_A\| \leq \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{T}} \ell_i(\mathbf{T})\| \leq \left[(\sqrt{t}B_x B_w + B_y)B_w B_x \beta B_{s'} \frac{B_x^2}{\sqrt{d}} \right] \|\mathbf{T}\|. \quad (60)$$

Thus we have proven equation 13 and next we prove equation 14.

• **Bounding the Laplacian of the Loss** For the second-order variations, we expand the Laplacian using the chain rule. By substituting both the Jacobian bound ($\beta B_{s'}$) and the Softmax Hessian bound ($\beta^2 B_{s''}$) from Lemma A.2 into $|\Delta_{\mathbf{W}_Q} \ell_i|$ and $|\Delta_{\mathbf{W}_K} \ell_i|$, we obtain the Laplacian bound.

Recalling, $\ell_i(\mathbf{T}) = \frac{1}{2} \|\mathbf{Y}_i - \hat{\mathbf{Y}}_i(\mathbf{T})\|_F^2 = \frac{1}{2} \text{Tr}(\mathbf{E}_i^\top \mathbf{E}_i)$, firstly we observe that,

$$\frac{\partial \ell_i}{\partial \mathbf{W}_{Q_{jk}}} = \frac{1}{2} \sum_{a,b} 2\mathbf{E}_{i,ab} \cdot \frac{\partial \mathbf{E}_{i,ab}}{\partial \mathbf{W}_{Q_{jk}}} = \sum_{a,b} \mathbf{E}_{i,ab} \cdot \frac{\partial (\hat{\mathbf{Y}}_{i,ab} - \mathbf{Y}_{i,ab})}{\partial \mathbf{W}_{Q_{jk}}}, \quad (61)$$

where $\mathbf{W}_{Q_{jk}}$ is the (j, k) -th element of the matrix \mathbf{W}_Q , which is a scalar.

Note that $\mathbf{Y}_{i,ab}$ is a constant and hence,

$$\begin{aligned} \frac{\partial^2 \ell_i}{\partial \mathbf{W}_{Q_{jk}}^2} &= \sum_{a,b} \left(\left[\frac{\partial \mathbf{E}_{i,ab}}{\partial \mathbf{W}_{Q_{jk}}} \right] \cdot \frac{\partial \hat{\mathbf{Y}}_{i,ab}}{\partial \mathbf{W}_{Q_{jk}}} + \mathbf{E}_{i,ab} \cdot \left[\frac{\partial^2 \hat{\mathbf{Y}}_{i,ab}}{\partial \mathbf{W}_{Q_{jk}}^2} \right] \right) \\ &= \sum_{a,b} \left(\frac{\partial \hat{\mathbf{Y}}_{i,ab}}{\partial \mathbf{W}_{Q_{jk}}} \right)^2 + \sum_{a,b} \left(\mathbf{E}_{i,ab} \cdot \frac{\partial^2 \hat{\mathbf{Y}}_{i,ab}}{\partial \mathbf{W}_{Q_{jk}}^2} \right) = \left\| \frac{\partial \hat{\mathbf{Y}}_i}{\partial \mathbf{W}_{Q_{jk}}} \right\|_F^2 + \text{Tr} \left(\mathbf{E}_i^\top \frac{\partial^2 \hat{\mathbf{Y}}_i}{\partial \mathbf{W}_{Q_{jk}}^2} \right). \end{aligned} \quad (62)$$

Towards analyzing the second term in the RHS above, consider the following derivative of $\mathbf{S}_i = \text{RowSoftMax}(\mathbf{M}_i)$.

$$\frac{\partial \mathbf{S}_{i,ab}}{\partial \mathbf{W}_{Q_{jk}}} = \sum_{c,d} \left(\frac{\partial \mathbf{S}_{i,ab}}{\partial \mathbf{M}_{i,cd}} \cdot \frac{\partial \mathbf{M}_{i,cd}}{\partial \mathbf{W}_{Q_{jk}}} \right) = \nabla \text{Softmax} \left[\frac{\partial \mathbf{M}_i}{\partial \mathbf{W}_{Q_{jk}}} \right]. \quad (63)$$

We have denoted the RHS to be $\nabla \text{Softmax} \left[\frac{\partial \mathbf{M}_i}{\partial \mathbf{W}_{Q_{jk}}} \right]$. This means the perturbation of the underlying weight $\mathbf{W}_{Q_{jk}}$ induces a directional change in the intermediate pre-activation matrix \mathbf{M}_i . By passing this direction through the Jacobian tensor of the Softmax operator and performing a tensor contraction, we compute the exact resulting variation in the output probability matrix \mathbf{S}_i .

To compute the second derivative $\frac{\partial^2 \mathbf{S}_i}{\partial \mathbf{W}_{Q_{jk}}^2}$, we differentiate the above result with respect to $\mathbf{W}_{Q_{jk}}$ again to get,

$$\frac{\partial^2 \mathbf{S}_{i,ab}}{\partial \mathbf{W}_{Q_{jk}}^2} = \sum_{c,d} \sum_{e,f} \left(\frac{\partial^2 \mathbf{S}_{i,ab}}{\partial \mathbf{M}_{i,cd} \partial \mathbf{M}_{i,ef}} \cdot \frac{\partial \mathbf{M}_{i,cd}}{\partial \mathbf{W}_{Q_{jk}}} \cdot \frac{\partial \mathbf{M}_{i,ef}}{\partial \mathbf{W}_{Q_{jk}}} \right) + \sum_{c,d} \left(\frac{\partial \mathbf{S}_{i,ab}}{\partial \mathbf{M}_{i,cd}} \cdot \frac{\partial^2 \mathbf{M}_{i,cd}}{\partial \mathbf{W}_{Q_{jk}}^2} \right). \quad (64)$$

As before, we denote the RHS in tensor notation as “ $\nabla^2 \text{Softmax} \left[\frac{\partial \mathbf{M}_i}{\partial \mathbf{W}_{Q_{jk}}}, \frac{\partial \mathbf{M}_i}{\partial \mathbf{W}_{Q_{jk}}} \right] + \nabla \text{Softmax} \left[\frac{\partial^2 \mathbf{M}_i}{\partial \mathbf{W}_{Q_{jk}}^2} \right]$ ”— the first term represents the second-order sensitivity of the Softmax operator (a 6th-order Hessian tensor) acting simultaneously as a

bilinear map on two identical first-order directional perturbations and the second term represents the first-order sensitivity (Jacobian tensor) acting on the second-order perturbation of the intermediate matrix M_i . Recalling $\hat{Y}_i(\mathbf{T}) = \mathbf{S}_i(\mathbf{T})\mathbf{X}_i\mathbf{W}_v$, and using the above notation we have

$$\begin{aligned}
 \frac{\partial^2 \hat{Y}_i}{\partial \mathbf{W}_{Q_{jk}}^2} &= \frac{\partial^2 \hat{\mathbf{S}}_i}{\partial \mathbf{W}_{Q_{jk}}^2} \mathbf{X}_i \mathbf{W}_v = \frac{\partial}{\partial \mathbf{W}_{Q_{jk}}} \frac{\partial \mathbf{S}_i}{\partial \mathbf{W}_{Q_{jk}}} \mathbf{X}_i \mathbf{W}_v = \frac{\partial}{\partial \mathbf{W}_{Q_{jk}}} \left(\nabla \text{Softmax} \left[\frac{\partial M_i}{\partial \mathbf{W}_{Q_{jk}}} \right] \right) \mathbf{X}_i \mathbf{W}_v \\
 &= \left(\nabla^2 \text{Softmax} \left[\frac{\partial M_i}{\partial \mathbf{W}_{Q_{jk}}}, \frac{\partial M_i}{\partial \mathbf{W}_{Q_{jk}}} \right] + \nabla \text{Softmax} \left[\frac{\partial^2 M_i}{\partial \mathbf{W}_{Q_{jk}}^2} \right] \right) \mathbf{X}_i \mathbf{W}_v.
 \end{aligned} \tag{65}$$

Note that $M_i = \frac{1}{\sqrt{d}} \mathbf{X}_i \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}_i^\top$ is linear in \mathbf{W}_Q , so $\frac{\partial^2 M_i}{\partial \mathbf{W}_{Q_{jk}}^2} \equiv 0$. From equation 64, using equation 50 and the Hessian bound from Lemma A.2, we have

$$\sum_{j,k} \left\| \frac{\partial^2 \mathbf{S}_i}{\partial \mathbf{W}_{Q_{jk}}^2} \right\|_F \leq \sum_{j,k} \|\nabla^2 \text{Softmax}(M_i)\|_F \cdot \left\| \frac{\partial M_i}{\partial \mathbf{W}_{Q_{jk}}} \right\|_F \cdot \left\| \frac{\partial M_i}{\partial \mathbf{W}_{Q_{jk}}} \right\|_F \leq \beta^2 B_{s''} \sum_{j,k} \left\| \frac{\partial M_i}{\partial \mathbf{W}_{Q_{jk}}} \right\|_F^2 \leq \beta^2 B_{s''} \frac{B_x^4}{d} \|\mathbf{W}_K\|_F^2. \tag{66}$$

So

$$|\Delta_{\mathbf{W}_Q} \ell_i| \leq (\beta B_{s'} B_x B_w)^2 \left[\frac{B_x^2}{\sqrt{d}} \right]^2 \|\mathbf{W}_K\|_F^2 + (\sqrt{t} B_x B_w + B_y) B_x B_w \beta^2 B_{s''} \frac{B_x^4}{d} \|\mathbf{W}_K\|_F^2. \tag{67}$$

By symmetry,

$$|\Delta_{\mathbf{W}_K} \ell_i| \leq (\beta B_{s'} B_x B_w)^2 \left[\frac{B_x^2}{\sqrt{d}} \right]^2 \|\mathbf{W}_Q\|_F^2 + (\sqrt{t} B_x B_w + B_y) B_x B_w \beta^2 B_{s''} \frac{B_x^4}{d} \|\mathbf{W}_Q\|_F^2. \tag{68}$$

We have,

$$\|\Delta_{\mathbf{T}} \hat{R}_A\| \leq \left\| \frac{1}{n} \sum_{i=1}^n \Delta_{\mathbf{T}} \ell_i(\mathbf{T}) \right\| \leq \|\Delta_{\mathbf{T}} \ell_i(\mathbf{T})\| \leq \frac{B_x^4}{d} \left((\beta B_{s'} B_x B_w)^2 + (\sqrt{t} B_x B_w + B_y) B_x B_w \beta^2 B_{s''} \right) \|\mathbf{T}\|^2. \tag{69}$$

□

C. Proof of Villani Conditions for Regression on Shallow Nets Under LoRA Constraints

This appendix provides the formal proof of Theorem 3.2, which establishes that the factor-regularized loss landscapes of shallow neural networks under LoRA constraints satisfy the Villani condition. By isolating the upper bound of the neural network's data fitting term, we can demonstrate how the factor regularization dominates as the parameters grow to infinity.

Lemma C.1. The corrected potential $V_\epsilon(\mathbf{T})$ satisfies the confining condition for all $\lambda > 0$.

Lemma C.2. For the i -th sample, the input is $\mathbf{x}_i \in \mathbb{R}^d$. Define $\mathbf{h}_i := \mathbf{V}^\top \mathbf{x}_i \in \mathbb{R}^r$, $\mathbf{s}_i := \mathbf{U} \mathbf{h}_i \in \mathbb{R}^p$, and

$$\mathbf{z}_i(\mathbf{U}, \mathbf{V}) := \mathbf{a}^\top \sigma(\mathbf{s}_i) = \sum_{j=1}^p a_j \sigma(s_{i,j}),$$

where σ acts component-wise. Defining $\mathbf{g}_i := \mathbf{a} \odot \sigma'(\mathbf{s}_i) \in \mathbb{R}^p$, it follows that the gradients of \mathbf{z}_i with respect to $\mathbf{U} \in \mathbb{R}^{p \times r}$ and $\mathbf{V} \in \mathbb{R}^{d \times r}$ are $\nabla_{\mathbf{U}} \mathbf{z}_i = \mathbf{g}_i \mathbf{h}_i^\top$ and $\nabla_{\mathbf{V}} \mathbf{z}_i = \mathbf{x}_i (\mathbf{U}^\top \mathbf{g}_i)^\top$, that is

$$\nabla_{\mathbf{U}} \mathbf{z}_i = [\mathbf{a} \odot \sigma'(UV^\top \mathbf{x}_i)] (\mathbf{V}^\top \mathbf{x}_i)^\top \quad \text{and} \quad \nabla_{\mathbf{V}} \mathbf{z}_i = \mathbf{x}_i [\mathbf{a} \odot \sigma'(UV^\top \mathbf{x}_i)]^\top \mathbf{U}. \tag{70}$$

Lemma C.3. The bound of $\ell'(\mathbf{z}_i) = (\mathbf{z}_i - \mathbf{y}_i)$ is $\sup_{i \in \{1, \dots, n\}} |\mathbf{z}_i - \mathbf{y}_i| = B_0 = \|\mathbf{a}\|_2 \sqrt{p} B_\sigma + B_y < \infty$.

Lemma C.4. The upper bound of $\|\nabla_{\mathbf{T}} \mathbf{z}_i\|$ is given by

$$\|\nabla_{\mathbf{T}} \mathbf{z}_i\| \leq B_{\sigma'} B_x \|\mathbf{a}\|_2 \|\mathbf{T}\| \quad \forall i = 1, 2, \dots, n, \tag{71}$$

and hence the upper bound of $\|\nabla_{\mathbf{T}} \mathcal{L}\|$ is given by

$$\|\nabla_{\mathbf{T}} \mathcal{L}\| \leq B_0 B_{\sigma'} B_x \|\mathbf{a}\|_2 \|\mathbf{T}\|. \tag{72}$$

Lemma C.5. The upper bound of $|\Delta_{\mathbf{T}}z_i|$ is given by

$$|\Delta_{\mathbf{T}}z_i| \leq B_{\sigma''} \|\mathbf{a}\|_1 B_x^2 \|\mathbf{T}\|^2 \quad \forall i = 1, 2, \dots, n, \quad (73)$$

and hence the upper bound of the absolute Laplacian $|\Delta_{\mathbf{T}}\mathcal{L}|$ is given by

$$|\Delta_{\mathbf{T}}\mathcal{L}| = \left| \frac{1}{n} \sum_{i=1}^n (\|\nabla_{\mathbf{T}}z_i\|^2 + (z_i - \mathbf{y}_i)\Delta_{\mathbf{T}}z_i) \right| \leq \left((B_{\sigma'} B_x \|\mathbf{a}\|_2)^2 + B_0 B_{\sigma''} \|\mathbf{a}\|_1 B_x^2 \right) \|\mathbf{T}\|^2. \quad (74)$$

The above lemmas are proved in Appendix D.

C.1. Proof of Theorem 3.2 for Loss in Definition 2.6

Proof. We analyze the potential defined in Definition 2.6:

$$\tilde{V}(\mathbf{T}) = \mathcal{L}(\mathbf{T}) + \frac{\lambda}{2} \|\mathbf{T}\|^2 \log(1 + \|\mathbf{T}\|^2), \quad (75)$$

where $\|\mathbf{T}\|$ refers to the 2-norm of \mathbf{T} read as a vector.

Let the regularization term be $R(\mathbf{T}) = \frac{\lambda}{2} \|\mathbf{T}\|^2 \log(1 + \|\mathbf{T}\|^2)$.

• **Analysis of the Gradient Term:** The gradient of the regularization term is:

$$\nabla R(\mathbf{T}) = \lambda \mathbf{T} \log(1 + \|\mathbf{T}\|^2) + \lambda \mathbf{T} \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2}. \quad (76)$$

Taking the norm and lower-bounding it, we get:

$$\|\nabla R(\mathbf{T})\| \geq \lambda \|\mathbf{T}\| \log(1 + \|\mathbf{T}\|^2). \quad (77)$$

By the gradient bound in Lemma C.4 (Equation 72), we have:

$$\begin{aligned} \|\nabla \tilde{V}(\mathbf{T})\|^2 &= (\|\nabla R(\mathbf{T})\| + \|\nabla \mathcal{L}(\mathbf{T})\|)^2 \\ &\geq \|\nabla R(\mathbf{T})\|^2 - 2\|\nabla R(\mathbf{T})\| \|\nabla \mathcal{L}(\mathbf{T})\| \\ &\geq (\|\nabla R(\mathbf{T})\|)^2 - 2(\|\nabla R(\mathbf{T})\|) \sup(\|\nabla \mathcal{L}(\mathbf{T})\|) \\ &= \lambda^2 \|\mathbf{T}\|^2 \log(1 + \|\mathbf{T}\|^2)^2 - 2B_0 B_{\sigma'} B_x \|\mathbf{a}\|_2 \|\mathbf{T}\|^2 \lambda \left(\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right). \end{aligned} \quad (78)$$

• **Analysis of the Laplacian Term:** To compute the Laplacian of the regularization term, we take the divergence of $\nabla R(\mathbf{T})$:

$$\Delta R(\mathbf{T}) = \sum_{k=1}^D \frac{\partial}{\partial T_k} \left(\lambda T_k \left[\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right] \right). \quad (79)$$

Applying the product rule yields:

$$\begin{aligned} \Delta R(\mathbf{T}) &= D\lambda \left[\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right] + \sum_{k=1}^D \lambda T_k \left[\frac{2T_k}{1 + \|\mathbf{T}\|^2} + \frac{2T_k}{(1 + \|\mathbf{T}\|^2)^2} \right] \\ &= D\lambda \left[\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right] + \frac{2\lambda \|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} + \frac{2\lambda \|\mathbf{T}\|^2}{(1 + \|\mathbf{T}\|^2)^2}. \end{aligned} \quad (80)$$

Observe that as $\|\mathbf{T}\| \rightarrow \infty$, $\Delta R(\mathbf{T}) = \mathcal{O}(\log \|\mathbf{T}\|^2)$. Applying Lemma C.5 (Equation 109), the Laplacian of the data term is bounded by a quadratic. Therefore, the total Laplacian is bounded by:

$$\Delta \tilde{V}(\mathbf{T}) \leq \left((B_{\sigma'} B_x \|\mathbf{a}\|_2)^2 + B_0 B_{\sigma''} \|\mathbf{a}\|_1 B_x^2 \right) \|\mathbf{T}\|^2 + \Delta R(\mathbf{T}). \quad (81)$$

• **Verifying the Villani Condition:** Substituting Equations 78 and 81 into the Villani limit expression for a given $s > 0$:

$$\begin{aligned}
 \lim_{\|\mathbf{T}\| \rightarrow \infty} \left(\frac{1}{s} \|\nabla \tilde{V}(\mathbf{T})\|^2 - \Delta \tilde{V}(\mathbf{T}) \right) &\geq \lim_{\|\mathbf{T}\| \rightarrow \infty} \left[\frac{1}{s} \left(\lambda^2 \|\mathbf{T}\|^2 \log(1 + \|\mathbf{T}\|^2) \right)^2 \right. \\
 &\quad \left. - 2B_0 B_{\sigma'} B_x \|\mathbf{a}\|_2 \lambda \|\mathbf{T}\|^2 \left(\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right) \right. \\
 &\quad \left. - \left((B_{\sigma'} B_x \|\mathbf{a}\|_2)^2 + B_0 B_{\sigma''} \|\mathbf{a}\|_1 B_x^2 \right) \|\mathbf{T}\|^2 + \Delta R(\mathbf{T}) \right] \\
 &= \lim_{\|\mathbf{T}\| \rightarrow \infty} \|\mathbf{T}\|^2 \left[\underbrace{\frac{\lambda^2}{s} (\log(1 + \|\mathbf{T}\|^2))^2}_{\text{Eq. 78} \rightarrow +\infty \text{ (Dominant)}} \right. \\
 &\quad \left. - \underbrace{\frac{2\lambda B_0 B_{\sigma'} B_x \|\mathbf{a}\|_2}{s} \left(\log(1 + \|\mathbf{T}\|^2) + \frac{\|\mathbf{T}\|^2}{1 + \|\mathbf{T}\|^2} \right)}_{\text{Eq. 78} \rightarrow +\infty \text{ (slower)}} \right. \\
 &\quad \left. - \underbrace{\left((B_{\sigma'} B_x \|\mathbf{a}\|_2)^2 + B_0 B_{\sigma''} \|\mathbf{a}\|_1 B_x^2 \right)}_{\text{Eq. 81 Constant}} - \underbrace{\frac{\Delta R(\mathbf{T})}{\|\mathbf{T}\|^2}}_{\text{Eq. 80} \rightarrow 0} \right]. \tag{82}
 \end{aligned}$$

The limit does diverge to $+\infty$:

$$\lim_{\|\mathbf{T}\| \rightarrow \infty} \left(\frac{1}{s} \|\nabla \tilde{V}(\mathbf{T})\|^2 - \Delta \tilde{V}(\mathbf{T}) \right) = \infty. \tag{83}$$

This demonstrates that the Villani condition is satisfied for all $\lambda > 0$ and $s > 0$. \square

C.2. Proof of Theorem 3.2 for Loss in Definition 2.7

Proof. We recall from Definition 2.7 that,

$$V_\epsilon(\mathbf{T}) = \mathcal{L}(\mathbf{T}) + \frac{\lambda}{2} (\|\mathbf{U}\|_F^{2+\epsilon} + \|\mathbf{V}\|_F^{2+\epsilon}).$$

By Lemma C.1 we know that V_ϵ is a confining function. Now, to show that $V_\epsilon(\mathbf{T})$ is Villani, we have to verify if the following is satisfied:

$$\frac{\|\nabla V_\epsilon(\mathbf{T})\|^2}{s} - \Delta V_\epsilon(\mathbf{T}) \rightarrow \infty \quad \text{as} \quad \|\mathbf{T}\| \rightarrow \infty. \tag{84}$$

• **Analysis of the Gradient Term:** Let's analyze the asymptotic behavior of the quantity $\mathcal{L}(\mathbf{T}) + \frac{\lambda}{2} (\|\mathbf{U}\|_F^{2+\epsilon} + \|\mathbf{V}\|_F^{2+\epsilon})$.

The gradient of the data term is,

$$\nabla \mathcal{L}(\mathbf{T}) = \frac{1}{n} \sum_{i=1}^n \ell'_i(\mathbf{z}_i) \nabla_{\mathbf{T}} (\mathbf{a}^\top \sigma(\mathbf{U}\mathbf{V}^\top \mathbf{x}_i)), \tag{85}$$

where $\ell'(z_i) = z_i - y_i$. The total gradient norm squared is $\|\nabla_{\mathbf{T}} \mathbf{z}_i\|^2 = \|\nabla_{\mathbf{U}} \mathbf{z}_i\|_F^2 + \|\nabla_{\mathbf{V}} \mathbf{z}_i\|_F^2$.

Since we have

$$\begin{aligned}
 \nabla \left[\frac{\lambda}{2} (\|\mathbf{U}\|_F^{2+\epsilon} + \|\mathbf{V}\|_F^{2+\epsilon}) \right] &= \left(\nabla_{\mathbf{U}} \left[\frac{\lambda}{2} (\|\mathbf{U}\|_F^2)^{\frac{2+\epsilon}{2}} \right], \nabla_{\mathbf{V}} \left[\frac{\lambda}{2} (\|\mathbf{V}\|_F^2)^{\frac{2+\epsilon}{2}} \right] \right) \\
 &= \frac{\lambda}{2} \left(\frac{2+\epsilon}{2} (\|\mathbf{U}\|_F^2)^{\frac{\epsilon}{2}} \cdot \nabla_{\mathbf{U}} (\|\mathbf{U}\|_F^2), \frac{2+\epsilon}{2} (\|\mathbf{V}\|_F^2)^{\frac{\epsilon}{2}} \cdot \nabla_{\mathbf{V}} (\|\mathbf{V}\|_F^2) \right), \\
 &= \frac{\lambda}{2} \left(\frac{2+\epsilon}{2} \|\mathbf{U}\|_F^\epsilon \cdot (2\mathbf{U}), \frac{2+\epsilon}{2} \|\mathbf{V}\|_F^\epsilon \cdot (2\mathbf{V}) \right), \\
 &= \frac{\lambda}{2} ((2+\epsilon) \|\mathbf{U}\|_F^\epsilon \mathbf{U}, (2+\epsilon) \|\mathbf{V}\|_F^\epsilon \mathbf{V})
 \end{aligned} \tag{86}$$

we get,

$$\|\nabla R_\varepsilon(\mathbf{T})\|^2 = \left(\frac{\lambda}{2}\right)^2 (2 + \varepsilon)^2 (\|\mathbf{U}\|_F^{2+2\varepsilon} + \|\mathbf{V}\|_F^{2+2\varepsilon}) \geq \left(\frac{\lambda}{2}\right)^2 (2 + \varepsilon)^2 2^{-\varepsilon} \|\mathbf{T}\|^{2+2\varepsilon}, \quad (87)$$

where we have used $\frac{(\|\mathbf{U}\|^2)^{1+\varepsilon} + (\|\mathbf{V}\|^2)^{1+\varepsilon}}{2} \geq \left(\frac{\|\mathbf{U}\|^2 + \|\mathbf{V}\|^2}{2}\right)^{1+\varepsilon} = 2^{-1-\varepsilon} \|\mathbf{T}\|^{2+2\varepsilon}$ by Jensen's Inequality.

Hence, substituting the gradient bound from Lemma C.4 (Equation 72), the gradient of the total potential $\nabla V_\varepsilon = \nabla \mathcal{L} + \lambda \mathbf{T}$ satisfies:

$$\begin{aligned} \|\nabla V_\varepsilon\|^2 &= \|\nabla \mathcal{L}\|^2 + \left\| \nabla \left[\frac{\lambda}{2} (\|\mathbf{U}\|_F^{2+\varepsilon} + \|\mathbf{V}\|_F^{2+\varepsilon}) \right] \right\|^2 \\ &\quad + \lambda \left\langle \frac{1}{n} \sum_{i=1}^n (\mathbf{z}_i - \mathbf{y}_i) \nabla_{\mathbf{T}} \mathbf{z}_i, ((2 + \varepsilon) \|\mathbf{U}\|_F^\varepsilon \mathbf{U}, (2 + \varepsilon) \|\mathbf{V}\|_F^\varepsilon \mathbf{V}) \right\rangle \\ &\geq \left(\frac{\lambda}{2}\right)^2 (2 + \varepsilon)^2 (\|\mathbf{U}\|_F^{2+2\varepsilon} + \|\mathbf{V}\|_F^{2+2\varepsilon}) \\ &\quad - \lambda \left\langle \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \mathbf{z}_i) \nabla_{\mathbf{T}} \mathbf{z}_i, ((2 + \varepsilon) \|\mathbf{U}\|_F^\varepsilon \mathbf{U}, (2 + \varepsilon) \|\mathbf{V}\|_F^\varepsilon \mathbf{V}) \right\rangle \\ &\geq \left(\frac{\lambda}{2}\right)^2 (2 + \varepsilon)^2 2^{-\varepsilon} \|\mathbf{T}\|^{2+2\varepsilon} - \lambda(2 + \varepsilon) B_0 B_{\sigma'} B_x \|\mathbf{a}\|_2 \|\mathbf{T}\|^{2+\varepsilon}. \end{aligned} \quad (88)$$

In the last inequality, we used $\|\mathbf{U}\|_F^{2+2\varepsilon} + \|\mathbf{V}\|_F^{2+2\varepsilon} \leq (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2)^{1+\varepsilon} = \|\mathbf{T}\|^{2+2\varepsilon}$.

• **Analysis of the Laplacian Term:** Since $\frac{\partial}{\partial U_{ij}} \|\mathbf{U}\|_F^{2+\varepsilon} = (2 + \varepsilon) \|\mathbf{U}\|_F^\varepsilon U_{ij}$ from the analysis above, we have

$$\frac{\partial^2}{\partial U_{ij}^2} \|\mathbf{U}\|_F^{2+\varepsilon} = (2 + \varepsilon) \varepsilon \|\mathbf{U}\|_F^{\varepsilon-2} U_{ij}^2 + (2 + \varepsilon) \|\mathbf{U}\|_F^\varepsilon. \quad (89)$$

So,

$$\Delta_{\mathbf{U}} \|\mathbf{U}\|_F^{2+\varepsilon} = (2 + \varepsilon) \|\mathbf{U}\|_F^\varepsilon \cdot \varepsilon + (2 + \varepsilon) \|\mathbf{U}\|_F^\varepsilon \cdot pr = (2 + \varepsilon)(\varepsilon + pr) \|\mathbf{U}\|_F^\varepsilon. \quad (90)$$

By symmetry,

$$\Delta_{\mathbf{V}} \|\mathbf{V}\|_F^{2+\varepsilon} = (2 + \varepsilon)(\varepsilon + dr) \|\mathbf{V}\|_F^\varepsilon. \quad (91)$$

Thus, we have:

$$\begin{aligned} \Delta R_\varepsilon(\mathbf{T}) &= \frac{\lambda}{2} [\Delta_{\mathbf{U}} \|\mathbf{U}\|_F^{2+\varepsilon} + \Delta_{\mathbf{V}} \|\mathbf{V}\|_F^{2+\varepsilon}] = \frac{\lambda}{2} (2 + \varepsilon) [(\varepsilon + pr) \|\mathbf{U}\|_F^\varepsilon + (\varepsilon + dr) \|\mathbf{V}\|_F^\varepsilon] \\ &\leq \frac{\lambda}{2} (2 + \varepsilon) [(\varepsilon + pr) \|\mathbf{T}\|^\varepsilon + (\varepsilon + dr) \|\mathbf{T}\|^\varepsilon] \\ &= \frac{\lambda}{2} (2 + \varepsilon) (2\varepsilon + D) \|\mathbf{T}\|^\varepsilon, \end{aligned} \quad (92)$$

where $D = (p + d)r$ is the total parameter dimension.

By Lemma C.5 (Equation 109), we have:

$$|\Delta_{\mathbf{T}} \mathcal{L}| \leq \left((B_{\sigma'} B_x \|\mathbf{a}\|_2)^2 + B_0 B_{\sigma''} \|\mathbf{a}\|_1 B_x^2 \right) \|\mathbf{T}\|^2. \quad (93)$$

Substituting Equation 88, Equation 92, and the data Laplacian bound into the Villani limit:

$$\begin{aligned}
 \lim_{\|\mathbf{T}\| \rightarrow \infty} \left(\frac{1}{s} \|\nabla V_\epsilon\|^2 - \Delta V_\epsilon \right) &\geq \lim_{\|\mathbf{T}\| \rightarrow \infty} \left[\frac{1}{s} \left(\left(\frac{\lambda}{2} \right)^2 (2+\epsilon)^2 2^{-\epsilon} \|\mathbf{T}\|^{2+2\epsilon} - \lambda(2+\epsilon) B_0 B_{\sigma'} B_x \|\mathbf{a}\|_2 2^{-\epsilon/2} \|\mathbf{T}\|^{2+\epsilon} \right) \right. \\
 &\quad \left. - \left((B_{\sigma'} B_x \|\mathbf{a}\|_2)^2 + B_0 B_{\sigma''} \|\mathbf{a}\|_1 B_x^2 \right) \|\mathbf{T}\|^2 + \frac{\lambda}{2} (2+\epsilon) (2\epsilon + D) \|\mathbf{T}\|^\epsilon \right] \\
 &= \lim_{\|\mathbf{T}\| \rightarrow \infty} \|\mathbf{T}\|^2 \left[\underbrace{\frac{1}{s} \left(\frac{\lambda}{2} \right)^2 (2+\epsilon)^2 2^{-\epsilon} \|\mathbf{T}\|^{2\epsilon}}_{\text{Eq. 88} \rightarrow +\infty} - \underbrace{\frac{\lambda}{s} (2+\epsilon) B_0 B_{\sigma'} B_x \|\mathbf{a}\|_2 2^{-\epsilon/2} \|\mathbf{T}\|^\epsilon}_{\text{Eq. 88} \rightarrow +\infty \text{ (slower)}} \right. \\
 &\quad \left. - \underbrace{\left((B_{\sigma'} B_x \|\mathbf{a}\|_2)^2 + B_0 B_{\sigma''} \|\mathbf{a}\|_1 B_x^2 \right)}_{\text{Eq. 109 Constant}} - \underbrace{\frac{\lambda}{2} (2+\epsilon) (2\epsilon + D) \|\mathbf{T}\|^{\epsilon-2}}_{\text{Eq. 92} \rightarrow 0} \right] \quad (94)
 \end{aligned}$$

$$= +\infty. \quad (95)$$

Since the quantity tends to positive infinity, the Villani condition is satisfied for all $\lambda > 0$ and $\epsilon > 0$. This completes the verification that the potential (Definition 2.7) induces the isoperimetric properties necessary for the Poincaré Inequality to hold. \square

D. Proofs of Intermediate Lemmas for Theorem 3.2

Proof of Lemma C.1. The Mean Square loss is non-negative, $\ell_i(\mathbf{W}) \geq 0$, thus $\mathcal{L}(\mathbf{T}) \geq 0$. The potential is bounded below by the factor regularization term:

$$V_\epsilon(\mathbf{T}) \geq \frac{\lambda}{2} \|\mathbf{T}\|^2. \quad (96)$$

Since $\lambda > 0$, the quadratic growth of $\frac{\lambda}{2} \|\mathbf{T}\|^2$ ensures $\lim_{\|\mathbf{T}\| \rightarrow \infty} V_\epsilon(\mathbf{T}) = +\infty$. This guarantees integrability of $e^{-\beta V_\epsilon}$ and normalizability of μ_β . \square

Proof of Lemma C.2. Since $d\mathbf{s}_i = d\mathbf{U} \mathbf{h}_i$, we have $d\mathbf{z}_i = (\mathbf{a} \odot \sigma'(\mathbf{s}_i))^\top d\mathbf{s}_i = \mathbf{g}_i^\top (d\mathbf{U} \mathbf{h}_i)$. This can be expressed in a form of inner product

$$\mathbf{g}_i^\top (d\mathbf{U} \mathbf{h}_i) = \text{tr}(\mathbf{h}_i \mathbf{g}_i^\top d\mathbf{U}) = \langle \mathbf{g}_i \mathbf{h}_i^\top, d\mathbf{U} \rangle,$$

so

$$\nabla_{\mathbf{U}} \mathbf{z}_i = \mathbf{g}_i \mathbf{h}_i^\top$$

Consider the derivative acting on $\mathbf{V}: d\mathbf{s}_i = \mathbf{U} d\mathbf{h}_i = \mathbf{U} (d\mathbf{V})^\top \mathbf{x}_i$, similarly

$$d\mathbf{z}_i = \mathbf{g}_i^\top \mathbf{U} (d\mathbf{V})^\top \mathbf{x}_i = \text{tr}((\mathbf{x}_i (\mathbf{U}^\top \mathbf{g}_i)^\top)^\top d\mathbf{V}) = \langle \mathbf{x}_i (\mathbf{U}^\top \mathbf{g}_i)^\top, d\mathbf{V} \rangle,$$

so

$$\nabla_{\mathbf{V}} \mathbf{z}_i = \mathbf{x}_i (\mathbf{U}^\top \mathbf{g}_i)^\top.$$

So the partial gradients w.r.t. the factors are:

$$\begin{aligned}
 \nabla_{\mathbf{U}} \mathbf{z}_i &= [\mathbf{a} \odot \sigma'(\mathbf{U} \mathbf{V}^\top \mathbf{x}_i)] (\mathbf{V}^\top \mathbf{x}_i)^\top \\
 \nabla_{\mathbf{V}} \mathbf{z}_i &= \mathbf{x}_i [\mathbf{a} \odot \sigma'(\mathbf{U} \mathbf{V}^\top \mathbf{x}_i)]^\top \mathbf{U}
 \end{aligned}$$

\square

Proof of Lemma C.3. Apply Cauchy-Schwarz, $|\mathbf{z}_i| \leq \|\mathbf{a}\|_2 \sqrt{p} B_\sigma$, recall $\sigma(\mathbf{s}_i) \in \mathbb{R}^p$. Since $|\mathbf{y}_i| \leq B_y$, we have

$$|\mathbf{z}_i - \mathbf{y}_i| \leq B_0 = \|\mathbf{a}\|_2 \sqrt{p} B_\sigma + B_y.$$

\square

1265 *Proof of Lemma C.4.* By Lemma C.2, $\|\nabla_{\mathbf{U}} \mathbf{z}_i\|_F = \|\mathbf{g}_i \mathbf{h}_i^\top\|_F = \|\mathbf{g}_i\|_2 \|\mathbf{h}_i\|_2$, where $\|\mathbf{g}_i\|_2 \leq B_{\sigma'} \|a\|_2$, $\|\mathbf{h}_i\|_2 = \|\mathbf{V}^\top \mathbf{x}_i\|_2 \leq$
 1266 $\|\mathbf{V}^\top\|_2 \|\mathbf{x}_i\|_2 = \|\mathbf{V}\|_2 \|\mathbf{x}_i\|_2 \leq \|\mathbf{V}\|_F B_x$. We have

$$1267 \quad \|\nabla_{\mathbf{U}} \mathbf{z}_i\|_F \leq (B_{\sigma'} \|a\|_2) (B_x \|\mathbf{V}\|_F) = B_{\sigma'} B_x \|a\|_2 \|\mathbf{V}\|_F. \quad (97)$$

1269 Similarly,

$$1270 \quad \|\nabla_{\mathbf{V}} \mathbf{z}_i\|_F \leq B_x \|\mathbf{U}\|_F (B_{\sigma'} \|a\|_2) = B_{\sigma'} B_x \|a\|_2 \|\mathbf{U}\|_F. \quad (98)$$

1273 Since, $d\mathbf{z}_i = \langle \nabla_{\mathbf{T}} \mathbf{z}_i, d\mathbf{T} \rangle$, combine 97 and 98 together to get

$$1274 \quad \|\nabla_{\mathbf{T}} \mathbf{z}_i\|^2 = \|\nabla_{\mathbf{U}} \mathbf{z}_i\|_F^2 + \|\nabla_{\mathbf{V}} \mathbf{z}_i\|_F^2 \leq (B_{\sigma'} B_x \|a\|_2)^2 (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2), \quad (99)$$

1277 we have $\|\nabla_{\mathbf{T}} \mathbf{z}_i\| \leq B_{\sigma'} B_x \|a\|_2 \|\mathbf{T}\|$ for $\forall i = 1, 2, \dots, n$. By Lemma C.3, we have,

$$1279 \quad \|\nabla_{\mathbf{T}} \mathcal{L}\| = \left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{z}_i - \mathbf{y}_i) \nabla_{\mathbf{T}} \mathbf{z}_i \right\| \leq \frac{1}{n} \sum_{i=1}^n |\mathbf{z}_i - \mathbf{y}_i| \|\nabla_{\mathbf{T}} \mathbf{z}_i\| \leq \frac{1}{n} \sum_{i=1}^n B_0 B_{\sigma'} B_x \|a\|_2 \|\mathbf{T}\| = \frac{1}{n} \cdot n \cdot B_0 B_{\sigma'} B_x \|a\|_2 \|\mathbf{T}\| \quad (100)$$

1282 \square

1284 *Proof of Lemma C.5.* We evaluate the data term Laplacian as follows:

$$1285 \quad \Delta \mathcal{L}(\mathbf{T}) = \frac{1}{n} \sum_{i=1}^n (\ell_i''(\mathbf{z}_i) \|\nabla_{\mathbf{T}} \mathbf{z}_i\|^2 + \ell_i'(\mathbf{z}_i) \Delta_{\mathbf{T}} \mathbf{z}_i) = \frac{1}{n} \sum_{i=1}^n (\|\nabla_{\mathbf{T}} \mathbf{z}_i\|^2 + (\mathbf{z}_i - \mathbf{y}_i) \Delta_{\mathbf{T}} \mathbf{z}_i). \quad (101)$$

1288 Recall, $\mathbf{h} = \mathbf{V}^\top \mathbf{x} \in \mathbb{R}^r$, $\mathbf{s} = \mathbf{U} \mathbf{h} \in \mathbb{R}^p$, and $z = \mathbf{a}^\top \sigma(\mathbf{s})$.

1290 For $\Delta_{\mathbf{T}} \mathbf{z}_i = \sum_{j=1}^p \sum_{k=1}^r \frac{\partial^2 \mathbf{z}_i}{\partial \mathbf{U}_{jk}^2} + \sum_{\ell=1}^d \sum_{k=1}^r \frac{\partial^2 \mathbf{z}_i}{\partial \mathbf{V}_{\ell k}^2}$, we can analyze term by term.

1292 Since, $\mathbf{s}_j = \sum_{k'} \mathbf{U}_{jk'} \mathbf{h}_{k'}$, we have $\frac{\partial \mathbf{s}_j}{\partial \mathbf{U}_{jk}} = \mathbf{h}_k$, so $\frac{\partial \mathbf{z}_i}{\partial \mathbf{U}_{jk}} = a_j \sigma'(\mathbf{s}_j) \frac{\partial \mathbf{s}_j}{\partial \mathbf{U}_{jk}} = a_j \sigma'(\mathbf{s}_j) \mathbf{h}_k$,

1294 and hence, $\frac{\partial^2 \mathbf{z}_i}{\partial \mathbf{U}_{jk}^2} = a_j \sigma''(\mathbf{s}_j) \left(\frac{\partial \mathbf{s}_j}{\partial \mathbf{U}_{jk}} \right)^2 = a_j \sigma''(\mathbf{s}_j) \mathbf{h}_k^2$.

1296 So,

$$1297 \quad \Delta_{\mathbf{U}} \mathbf{z}_i = \sum_{j=1}^p \sum_{k=1}^r a_j \sigma''(\mathbf{s}_j) \mathbf{h}_k^2 = \left(\sum_{k=1}^r \mathbf{h}_k^2 \right) \left(\sum_{j=1}^p a_j \sigma''(\mathbf{s}_j) \right) = \|\mathbf{h}\|_2^2 \sum_{j=1}^p a_j \sigma''(\mathbf{s}_j). \quad (102)$$

1301 Similarly, for $\mathbf{h}_{k'} = \sum_{\ell'=1}^d \mathbf{V}_{\ell' k'} \mathbf{x}_{\ell'}$, we can express s with

$$1302 \quad \mathbf{s}_j = \sum_{k'=1}^r \mathbf{U}_{jk'} \mathbf{h}_{k'} = \sum_{k'=1}^r \mathbf{U}_{jk'} \left(\sum_{\ell'=1}^d \mathbf{V}_{\ell' k'} \mathbf{x}_{\ell'} \right),$$

1306 so $\frac{\partial \mathbf{s}_j}{\partial \mathbf{V}_{\ell k}} = \mathbf{U}_{jk} \mathbf{x}_{\ell}$, hence $\frac{\partial \mathbf{z}_i}{\partial \mathbf{V}_{\ell k}} = \sum_{j=1}^p a_j \sigma'(\mathbf{s}_j) \frac{\partial \mathbf{s}_j}{\partial \mathbf{V}_{\ell k}} = \sum_{j=1}^p a_j \sigma'(\mathbf{s}_j) \mathbf{U}_{jk} \mathbf{x}_{\ell}$ and $\frac{\partial^2 \mathbf{z}_i}{\partial \mathbf{V}_{\ell k}^2} = \sum_{j=1}^p a_j \left[\sigma''(\mathbf{s}_j) \frac{\partial \mathbf{s}_j}{\partial \mathbf{V}_{\ell k}} \right] \mathbf{U}_{jk} \mathbf{x}_{\ell}$.

1308 So

$$1309 \quad \Delta_{\mathbf{V}} z = \sum_{\ell=1}^d \sum_{k=1}^r \frac{\partial^2 z}{\partial \mathbf{V}_{\ell k}^2} = \sum_{\ell=1}^d \sum_{k=1}^r \left(\sum_{j=1}^p a_j \sigma''(\mathbf{s}_j) \mathbf{U}_{jk}^2 \mathbf{x}_{\ell}^2 \right) = \sum_{j=1}^p a_j \sigma''(\mathbf{s}_j) \left(\sum_{\ell=1}^d \mathbf{x}_{\ell}^2 \right) \left(\sum_{k=1}^r \mathbf{U}_{jk}^2 \right), \quad (103)$$

1313 where $\mathbf{U}_j = \sum_{k=1}^r \mathbf{U}_{jk}^2$.

1314 So $\Delta_{\mathbf{T}} \mathbf{z}_i = \Delta_{\mathbf{U}} \mathbf{z}_i + \Delta_{\mathbf{V}} \mathbf{z}_i = \|\mathbf{h}\|_2^2 \sum_{j=1}^p a_j \sigma''(\mathbf{s}_j) + \|\mathbf{x}\|_2^2 \sum_{j=1}^p a_j \sigma''(\mathbf{s}_j) \|\mathbf{U}_j\|_2^2$.

1315 Determine the bound,

$$1316 \quad |\Delta_{\mathbf{U}} \mathbf{z}_i| = \left| \|\mathbf{h}\|_2^2 \sum_j a_j \sigma''(\mathbf{s}_j) \right| \leq \|\mathbf{h}\|_2^2 \sum_j |a_j| |\sigma''(\mathbf{s}_j)| \leq \|\mathbf{h}\|_2^2 \|a\|_1 B_{\sigma''}. \quad (104)$$

Since $\|\mathbf{h}\|_2 = \|\mathbf{V}^\top \mathbf{x}\|_2 \leq \|\mathbf{V}\|_2 \|\mathbf{x}\|_2 \leq \|\mathbf{V}\|_F B_x$, We have

$$|\Delta_{\mathbf{U}} z| \leq B_{\sigma''} \|a\|_1 B_x^2 \|\mathbf{V}\|_F^2. \quad (105)$$

Similarly,

$$\begin{aligned} |\Delta_{\mathbf{V}} z| &= \left| \|\mathbf{x}\|_2^2 \sum_j a_j \sigma''(\mathbf{s}_j) \|\mathbf{U}_j\|_2^2 \right| \leq \|\mathbf{x}\|_2^2 \sum_j |a_j| |\sigma''(\mathbf{s}_j)| \|\mathbf{U}_j\|_2^2 \leq B_x^2 B_{\sigma''} \sum_j |a_j| \|\mathbf{U}_j\|_2^2 \\ &\leq \|a\|_1 \sum_j \|\mathbf{U}_j\|_2^2 = \|a\|_1 \|\mathbf{U}\|_F^2. \end{aligned} \quad (106)$$

Combine equation 105 and 106, we have,

$$|\Delta_T z| \leq B_{\sigma''} \|a\|_1 B_x^2 (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) = B_{\sigma''} \|a\|_1 B_x^2 \|\mathbf{T}\|^2. \quad (107)$$

So for $(z_i - \mathbf{y}_i) \Delta_T z_i$, By Lemma C.3:

$$|(z - \mathbf{y}) \Delta_T z| \leq |z - \mathbf{y}| |\Delta_T z| \leq B_0 B_{\sigma''} \|a\|_1 B_x^2 \|\mathbf{T}\|^2. \quad (108)$$

Combine equation 99 and 108,

$$\Delta_T \mathcal{L} = \frac{1}{n} \sum_{i=1}^n (\|\nabla_{\mathbf{T}} z_i\|^2 + (z_i - \mathbf{y}_i) \Delta_T z_i) \leq ((B_{\sigma'} B_x \|a\|_2)^2 + B_0 B_{\sigma''} \|a\|_1 B_x^2) \|\mathbf{T}\|^2. \quad (109)$$

□

E. Proofs of SDE convergence

Proof of Theorem 3.3. Let $\tilde{L}^{(k)}(\mathbf{T})$ denote any of the four regularized potentials defined in Definition 2.4 (\tilde{V}_{ATT}), Definition 2.5 ($V_{\epsilon, \text{ATT}}$), Definition 2.6 (V), and Definition 2.7 (V_ϵ), where $k \in \{1, 2, 3, 4\}$ indexes the specific model and regularization choice.

For any $k \in \{1, 2, 3, 4\}$, perform the following decomposition of the excess risk:

$$\mathbb{E}[\tilde{L}^{(k)}(\mathbf{T}_t)] - \tilde{L}^{(k)*} = \underbrace{\left(\mathbb{E}[\tilde{L}^{(k)}(X_s^{(k)}(\infty))] - \tilde{L}^{(k)*} \right)}_{\varepsilon^{(k)}(s)} + \underbrace{\left(\mathbb{E}[\tilde{L}^{(k)}(\mathbf{T}_t)] - \mathbb{E}[\tilde{L}^{(k)}(X_s^{(k)}(\infty))] \right)}_{\leq D^{(k)}(s, p_0) e^{-\lambda_s^{(k)} t}} \quad (110)$$

By Proposition 5 of (Shi et al., 2023),

$$\varepsilon^{(k)}(s) \leq A^{(k)} s \quad (111)$$

and Proposition 4 of (Shi et al., 2023),

$$\mathbb{E}[\tilde{L}^{(k)}(\mathbf{T}_t)] - \mathbb{E}[\tilde{L}^{(k)}(X_s^{(k)}(\infty))] \leq C^{(k)}(s) \|p_0 - \mu_s^{(k)}\|_{L^2((\mu_s^{(k)})^{-1})} e^{-\lambda_s^{(k)} t}, \quad (112)$$

where $p_0 \in L^2((\mu_s^{(k)})^{-1})$ is the initial probability density of the SDE (1).

Since $s \leq \min \left\{ \frac{\epsilon}{2A^{(k)}}, S^{(k)} \right\}$, and the time horizon $t \geq \frac{1}{\lambda_s^{(k)}} \log \left(\frac{2D^{(k)}(s, p_0)}{\epsilon} \right)$, by Corollary 6 of (Shi et al., 2023) we have

$$\mathbb{E}[\tilde{L}^{(k)}(\mathbf{T}_t)] - \tilde{L}^{(k)*} \leq \epsilon. \quad (113)$$

□