
Trainable Transformer in Transformer

Abhishek Panigrahi Sadhika Malladi Mengzhou Xia Sanjeev Arora

Department of Computer Science

Princeton University

{ap34, smalladi, mengzhou, arora}@cs.princeton.edu

Abstract

Recent works attribute the capability of in-context learning (ICL) in large pre-trained language models to implicitly simulating and fine-tuning an internal model (e.g., linear or 2-layer MLP) during inference. However, such constructions require large memory overhead, which makes simulation of more sophisticated internal models intractable. In this work, we propose a new efficient construction, *Transformer in Transformer* (in short, TINT), that allows a transformer to simulate and fine-tune more complex models during inference (e.g., pre-trained language models). In particular, we introduce innovative approximation techniques that allow a TINT model with less than 2 billion parameters to simulate and fine-tune a 125 million parameter transformer model within a single forward pass. TINT accommodates many common transformer variants and its design ideas also improve the efficiency of past instantiations of simple models inside transformers. We conduct end-to-end experiments to validate the internal fine-tuning procedure of TINT on various language modeling and downstream tasks. For example, even with a limited one-step budget, we observe TINT for a OPT-125M model improves performance by 4 – 16% absolute on average compared to OPT-125M. These findings suggest that large pre-trained language models are capable of performing intricate subroutines. To facilitate further work, a modular and extensible codebase for TINT is open-sourced.¹

1 Introduction

Transformers [42] have brought about a revolution in language modeling, and scaling model size has enabled significant advancements in capabilities [6, 8]. One such capability [46] is in-context learning (ICL), where language models “learn” from given training exemplars in the context and subsequently predict the label of a test example within a single inference pass.

Several works [1, 11, 43] propose that ICL occurs when the large (“simulator”) model mimics and trains a smaller and simpler auxiliary model —such as a linear or 2-layer MLP model— on the in-context data. A crucial limitation of previous works is the large number of parameters needed for a simulator to perform such a complex subroutine during its forward pass, which restricts the simulator to performing very few training steps on fairly simple models. For example, simulating training of a linear layer can require tens of millions of parameters [1], and extending the simulator to train a larger model would require a simulator with trillions of parameters.

The current work shows that minor modifications to the standard transformer architecture allow it to efficiently simulate and approximately train an internal *auxiliary* transformer during a single inference pass (Section 2). We call our architecture *Transformer in Transformer*, or TINT in short. We show how TINT can internally simulate and train several popular and capable Transformer models such as GPT [33], OPT [53], and other variants [41, 36]. TINT requires fewer than two billion parameters to

¹https://github.com/abhishekpanigrahi1996/transformer_in_transformer

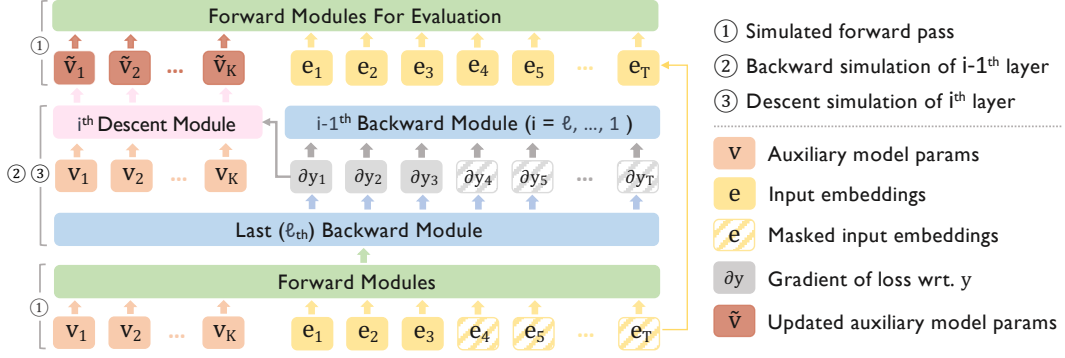


Figure 1: The overall structure of TINT. Each Forward, Backward, and Descent module is represented using combinations of linear, self-attention, layernorm, and activation layers. The input consists of prefix embeddings, that represent relevant auxiliary model parameters in each layer, input token embeddings, and a binary prefix mask to separate the train and evaluation segments of the input. Auxiliary model parameters are updated in the descent module using the training part, and the updated prefix tokens are transferred to the forward modules via residual connections for evaluating the rest.

internally simulate and update an auxiliary transformer with 125 million parameters (e.g., GPT-2 or OPT-125M). The scale of our construction is crucial to its significance, as it suggests that even transformers of moderate scale can explicitly learn from context during inference.

We validate our approach with end-to-end experiments on many language modeling and downstream tasks. Results demonstrate that a TINT model constructed to simulate and tune an OPT-125M model leads to a perplexity reduction of 0.3 to 0.7 points in language modeling. Additionally, TINT learns from in-context exemplars in the few-shot setting, resulting in an absolute gain of 12% to 16% over the auxiliary model. TINT can also learn from the context tokens of the evaluation inputs in the zero-shot setting, leading to an absolute performance improvement of up to 4% when no explicit exemplars are provided (Section 3). To the best of our knowledge, TINT is the first simulator to undergo such a comprehensive end-to-end evaluation on standard language tasks. In contrast, previous studies primarily conducted probing tests on transformers pre-trained using synthetic datasets or lacked empirical validation [1, 43, 14], likely due to the immense scale required by their constructions.

To summarize, our work improves over prior works in three crucial ways.

1. **Expressiveness:** We allow the auxiliary model to be a complex, general-purpose transformer, which requires significant technical innovation beyond past work using linear and MLP auxiliary models. TINT can internally simulate and train popular moderate- and large-scale transformers, including OPT [53], LLaMA [41], BLOOM [36], and Pythia [5].
2. **Efficiency:** We focus on making TINT orders of magnitude smaller, reducing the construction size from trillions of parameters to a few billion, despite simulating much complex auxiliary models. This efficiency is driven by a number of approximations (see Section 2.2).
3. **Empirical Validation:** TINT can match the performance of fine-tuning the auxiliary model explicitly, thereby validating the approximations used in the construction. In addition, we provide an extensible codebase for further training and evaluation of our proposed architecture (Section 3).

2 Our Construction

2.1 Overview

We design a transformer architecture, dubbed TINT, to perform the forward, backward, and parameter update operations of a smaller so-called *auxiliary* model over the course of a single inference pass (Figure 1). The first few layers of TINT simulate the forward pass, and then backpropagation and gradient updates are performed in parallel in the next several layers of the model. The final layers of TINT simulate another forward pass through the auxiliary model with the updated parameters, and TINT directly outputs the result. This procedure requires that TINT has (1) read and write access to the auxiliary model weights, and (2) read access to labelled training data.

For (1), we read from and write to the token embeddings of the first few tokens in the input (i.e., *prefix embeddings*, see Definition 2.1 and Section 2.3). For (2), we designate the first few input tokens as training data. If the input contains in-context demonstrations, then TINT performs a natural operation: training on in-context demonstrations and testing on the last example. However, the input can also be standard text tokens without additional formatting or labelling, in which case TINT performs gradient descent to adapt to the first few tokens in a context before being evaluated on the rest of the context. This procedure is known as dynamic evaluation, which we describe in further detail in Appendix L.

2.2 Key Components

Due to space constraints, we defer a complete description of the TINT to the appendix. Here, we detail the modifications and approximations introduced in our architecture to enable efficient simulation of training an internal auxiliary model. Experiments in Section 3 verify that these approximations do not harm performance. Numbers in parentheses indicate the parameter saving factor when relevant.

1. **Prefix embeddings** ($5\times$ compared to Wei et al. [45], Perez et al. [30]): As described in Section 2.3, we use the token embeddings of the first few inputs (i.e., the *prefix*, see Definition 2.1) to represent the relevant auxiliary model weights at each layer.
2. **H_{sim} -split linear operations** ($H_{\text{sim}}\times$): We parallelize expensive linear operations by splitting them across attention heads (Section 2.4).
3. **Linear attention**: We use linear attention modules to perform the forward, backward, and gradient operations for an auxiliary model linear layer. Softmax attention also suffices but requires more parameters and incurs an approximation error (Theorem 2.5). We use softmax attention modules to simulate the auxiliary model attention modules in TINT.
4. **First order gradients** ($4\times$): We use the first-order term of the gradient for layer normalization and activation functions (Appendix A.1).
5. **Only train the value vectors of the attention** ($5\times$): We only update and backpropagate through the value vectors of the attention layers. TINT is designed to simulate only a few steps of training the auxiliary model, so we expect this approximation not to drastically modify the final performance. Our experiments in Tables 1 and 2 validate that it does not hurt performance. Moreover, in Theorem A.8, we formally show that under certain conditions, backpropagating through just the value vectors can be arbitrarily accurate to standard backpropagation.
6. **Parameter sharing** ($3\times$ or $4\times$): We save $3\times$ parameters by applying the same forward module in TINT to simulate the query, key, and value computation of the auxiliary model’s self-attention module (Appendix A.3). Similarly, we divide the feedforward layer in the auxiliary model into 4 sub-layers and save $4\times$ parameters by employing a single TINT module to simulate each sub-layer.

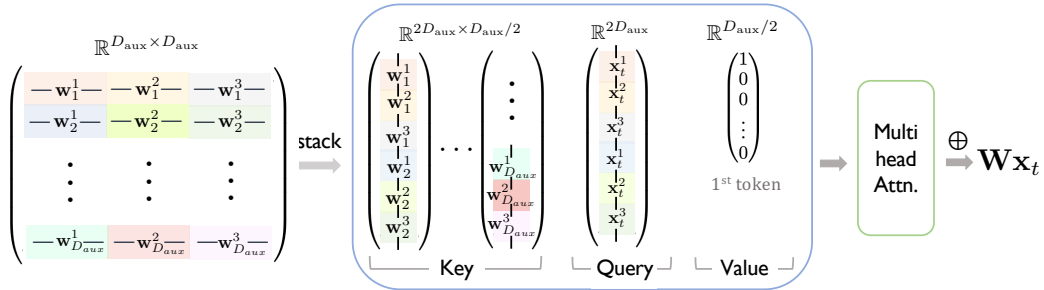


Figure 2: TINT simulates the forward pass of a linear layer as a H -head ($H = 6$ here) attention layer, with parameters of the auxiliary model as the key, the encodings of input tokens as the query, and the positional one-hot vector of the prefix embeddings as the value. We omitted the identical transformation for key, query, and value matrices for simplicity.

Notation: Let D denote the embedding dimension for a token and T denote the length of an input sequence. H denotes the number of attention heads. With the exception of contextual embeddings,

we use subscripts to indicate if the quantity is from TINT or from the auxiliary model. For example, D_{aux} refers to the embedding dimension and D_{sim} refers to the TINT embedding dimension. For contextual embeddings, we use $e_t^{(\ell)} \in \mathbb{R}^{D_{\text{sim}}}$ to denote activations in TINT and $x_t^{(\ell)} \in \mathbb{R}^{D_{\text{aux}}}$ to denote activations in the auxiliary model, where ℓ is the layer and t is the sequence position. When convenient, we drop the superscript that represents the layer index and the subscript that represents the position index. For a matrix A , a_j refers to its j th row, and for any vector b , b_j refers to its j th element. TINT uses one-hot positional embeddings $\{p_i^{\text{TINT}} \in \mathbb{R}^{T_{\text{sim}}}\}_{i \leq T_{\text{sim}}}$. For illustration, we ignore the bias parameters here but discuss them later in the appendix.

2.3 Operating on an auxiliary model with prefix embeddings

The straightforward way to simulate the forward pass of the auxiliary model would be to store its weights in the simulator’s weights and run a forward pass as usual. However, this gives the simulator no way to update the weights of the auxiliary model, since the simulator cannot modify its own weights during a forward pass. The only way to update the auxiliary model weights is by storing them in model *activations* that can be accessed and modified over the course of a forward pass.

Wei et al. [45], Perez et al. [30] model the simulator after a Turing machine, where each $e_t^{(\ell)} \in \mathbb{R}^{D_{\text{sim}}}$ acts as a workspace for operations, and data is copied between workspaces and memory using attention operations. Memory can either be allocated in a token embedding, thereby increasing the embedding size D_{sim} [1], or passed into the model as additional tokens, thereby increasing the simulator’s input sequence length T_{sim} . Both strategies increase the size of the construction, and using attention modules for copy operations results in a drastic scaling. For example, if $D_{\text{aux}} = 768$, a dot product with weight $w \in \mathbb{R}^{768}$, i.e. $\langle w, x_t^{(\ell)} \rangle$, requires at least 8.7 million parameters in the simulator².

Alternatively, storing memory as context tokens and allowing the attention modules to attend to those tokens removes the need for copying operations [14]. Then, a dot product with weight $w \in \mathbb{R}^{768}$, i.e. $\langle w, x_t^{(\ell)} \rangle$, requires 1.7 million parameters only. Naive implementation is problematic since the TINT attention module grows quadratically with the sequence length T_{sim} . We thus define *prefix embeddings* in TINT to contain only the relevant auxiliary parameters at each layer.

Definition 2.1 (Prefix Embeddings). $\{v_j^{(\ell)}\}_{j=1}^K$ denotes the K prefix embeddings at the ℓ th layer in TINT. These contain relevant auxiliary model weights or simulated activations.

Prefix embeddings permit efficient parallelization across attention heads while keeping the embedding dimension D_{sim} small, as we demonstrate below.

2.4 Stacking in prefix-tokens, H_{sim} -split linear operations and Linear attention

We motivate three parameter-efficient techniques using a $D_{\text{aux}} \times D_{\text{aux}}$ linear layer as a case study. The linear layer is applied token-wise, so we consider a single position t without loss of generality.

Definition 2.2 (Linear layer). For a weight $W \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$, a linear layer takes $x \in \mathbb{R}^{D_{\text{aux}}}$ as input and outputs $y = Wx$.

Stacking: We compute $\langle w_i, x_t \rangle$ for all $i \in [D_{\text{aux}}]$, where w_i denotes the i th row of W . To do so, the TINT input embedding e_t must contain x_t in its first D_{aux} coordinates, and the weights $\{w_i\}$ are in prefix embeddings $\{v_j\}$ (Definition 2.1). A first attempt is to put each w_i in its own v_i vector, which requires $K = D_{\text{aux}}$ prefix embeddings at the start of the sequence. For GPT-2, $D_{\text{aux}} = 768$, so this strategy will not allow the TINT to accept many standard language context tokens; moreover, the attention modules in the TINT will grow quadratically with input length. Therefore, we stack S weights on top of each other to form each prefix embedding v_i . S drives a trade-off between the embedding dimension of the TINT, $D_{\text{sim}} := D_{\text{aux}}S$, and the context length to the TINT, $T_{\text{sim}} := K + T_{\text{aux}}$. We set $S = 4$.

Attention Module: We use a self-attention module in TINT to compute the matrix-vector product. We modify the usual attention layer to also include the one-hot position embeddings $\{p_i^{\text{TINT}} \in \mathbb{R}^{T_{\text{sim}}}\}_{i \leq T_{\text{sim}}}$. Here, we use a single attention head (see Definition C.1 for multi-head attention).

²The copy attention will require 1.7 million parameters, while the dot product with a feedforward module (following [1]) will require > 7 million parameters.

Definition 2.3 (TINT self-attention with single head). For parameters $\{\mathbf{W}_Q^{\text{TINT}}, \mathbf{W}_K^{\text{TINT}}, \mathbf{W}_V^{\text{TINT}} \in \mathbb{R}^{D_{\text{sim}} \times D_{\text{sim}}}\}$, $\{\mathbf{W}_Q^p, \mathbf{W}_K^p, \mathbf{W}_V^p \in \mathbb{R}^{D_{\text{sim}} \times T_{\text{sim}}}\}$, the self-attention layer with single attention head and a function $f_{\text{attn}} : \mathbb{R}^{T_{\text{sim}}} \rightarrow \mathbb{R}^{T_{\text{sim}}}$ takes a sequence $\{\hat{\mathbf{e}}_t \in \mathbb{R}^{D_{\text{sim}}}\}_{t \leq T_{\text{sim}}}$ as input and outputs $\{\tilde{\mathbf{e}}_t \in \mathbb{R}^{D_{\text{sim}}}\}_{t \leq T_{\text{sim}}}$, such that

$$\begin{aligned} \tilde{\mathbf{e}}_t &= \sum_{j \leq T_{\text{sim}}} a_{t,j} \mathbf{v}_j, \quad \text{where } a_{t,j} = f_{\text{attn}}(\mathbf{K} \mathbf{q}_t)_j, \\ \mathbf{q}_t &= \mathbf{W}_Q^{\text{TINT}} \hat{\mathbf{e}}_t + \mathbf{W}_Q^p \mathbf{p}_t, \quad \mathbf{k}_t = \mathbf{W}_K^{\text{TINT}} \hat{\mathbf{e}}_t + \mathbf{W}_K^p \mathbf{p}_t, \quad \mathbf{v}_t = \mathbf{W}_V^{\text{TINT}} \hat{\mathbf{e}}_t + \mathbf{W}_V^p \mathbf{p}_t \text{ for all } t \leq T_{\text{sim}}, \\ \text{and } \mathbf{K} &\in \mathbb{R}^{T_{\text{sim}} \times D_{\text{sim}}} \text{ is the key matrix defined with its rows as } \{\mathbf{k}_t\}_{t \leq T_{\text{sim}}}. \end{aligned}$$

$\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t$ are the query, key, and value vectors at position t , and $a_{t,j}$ is the attention score between tokens at position t and j . f_{attn} can be either linear or softmax functions, and the corresponding layers are referred to as linear and softmax self-attention respectively.

Remark 2.4. In order to compute and backpropagate the loss during inference, the self-attention layers in TINT need to be non-causal on the first few tokens of the input.³ Explicit masks apply bidirectional attention to the input to backpropagate auxiliary self-attention layers (Appendix A.2).

TINT Linear Forward module (Figure 2): A first attempt would be to use different attention heads to operate on different rows; however, this uses S attention heads whereas large transformers usually have many more heads. Moreover, the output from the multi-head attention would need to be reorganized by a $D_{\text{sim}} \times D_{\text{sim}}$ linear layer before it could be reduced efficiently via summation. We instead parallelize across more attention heads to ensure the resulting output can easily be compiled: crucially, we shard each individual weight into S' parts. We set S and S' such that $H_{\text{sim}} = S \times S'$, using all available heads to parallelize the dot products.

H_{sim} -split linear operations (Appendix D.1): The output resulting from the attention module has shape $(D_{\text{sim}}/H_{\text{sim}}) \times H_{\text{sim}}$ and is sparse. To complete linear forward pass, we need to sum and aggregate the appropriate terms to form a D_{sim} -length vector with $\mathbf{W}\mathbf{x}$ in the first D_{aux} coordinates. Straightforwardly summing along an axis aggregates incorrect terms, since the model was sharded. Rearranging the entire matrix to enable easy summation requires an additional $D_{\text{sim}} \times D_{\text{sim}}$ linear layer. But we can save a $H_{\text{sim}} \times$ parameters by leveraging the local structure of the attention output. We space out the results across the $D_{\text{sim}}/H_{\text{sim}}$ rows and then sum along the H_{sim} columns to get the desired D_{sim} -length vector. This requires $D_{\text{sim}}^2/H_{\text{sim}} + D_{\text{sim}}H_{\text{sim}}$ parameters. Note that leveraging local structure also compresses the constructions for the TINT's backpropagation modules of layer normalization and activation operations (Appendices F and G).

Linear Attention: The above construction uses a linear attention mechanism instead of the canonical softmax. Here, we show that any linear attention module with bounded entries can be approximated by softmax attention with a few additional parameters. Thus, we often use linear attention in TINT.

Theorem 2.5 (Informal, c.f. Theorem C.2). *For any $\epsilon > 0$, $B > 0$, and a linear attention module with H heads and bounded parameters, there exists a softmax attention module with $2H_{\text{sim}}$ attention heads and $4 \times$ additional parameters, such that on every sequence of inputs with B -bounded norm, the output sequences of the softmax attention and the linear attention differ by $\mathcal{O}(\epsilon)$ at each position.*

3 Verification of TINT

We conduct experiments on TINT constructed using GPT2 and OPT-125M as auxiliary models to validate that the approximations introduced (Section 2.2) do not significantly harm the capability of TINT to simulate and train an internal model.

In language modeling (Table 2), the perplexity decreases using TINT, especially as the training proportion increases. For downstream tasks (Table 1), explicit internal training within TINT surpasses vanilla zero-shot evaluation and in-context learning, even with a limited budget of a single forward pass. Moreover, TINT achieves a performance comparable to dynamic evaluation, indicating that the approximations made during its construction largely preserve its effectiveness for fine-tuning. Though calibration may not always be beneficial in every setting,⁴ we observe that the efficacy of

³Similar prefix models have been developed in [34, 24].

⁴Such inconsistencies in the calibration method have been observed in previous works [6].

Table 1: Zero-shot and few-shot in-context learning results across 7 downstream tasks. All the few-shot results are averaged over three training seeds. TINT consistently surpasses its auxiliary model and achieves comparable performance to dynamic evaluation. TINT outperforms auxiliary models by 3 – 4% and 12 – 16% absolute points on average in 0-shot and 32-shot experiments respectively. TINT performs competitively with a similar-sized pre-trained model (OPT-1.3B) in both 0-shot and 32-shot settings. We show the standard deviation for few-shot settings in parentheses.

Model	Shots	Subj	AGNews	SST2	CR	MR	MPQA	Amazon	Avg.
<i>Without Calibration</i>									
OPT-125M	0	64.0	66.0	70.5	64.5	71.0	68.0	76.5	68.6
OPT-1.3B	0	59.0	55.5	54.0	50.5	52.5	74.0	57.0	57.5
OPT-125M DYNA. EVAL	0	71.0	67.0	79.5	71.5	70.0	68.0	85.5	73.2
OPT-125M TINT	0	67.5	66.0	76.5	69.0	76.0	70.5	78.5	72.0
OPT-125M	32	58.7 _(4.9)	33.7 _(8.4)	50.8 _(1.2)	51.3 _(1.9)	50.0 _(0.0)	54.3 _(2.5)	55.0 _(6.7)	50.5 _(1.9)
OPT-1.3B	32	74.2 _(6.1)	71.3 _(5.3)	89.8 _(3.6)	71.5 _(4.5)	68.3 _(6.1)	81.7 _(3.3)	70.3 _(9.9)	75.3 _(0.4)
OPT-125M DYNA. EVAL	32	78.0 _(1.4)	66.7 _(1.6)	71.5 _(1.4)	73.7 _(3.3)	72.0 _(0.0)	80.7 _(0.6)	79.8 _(0.2)	74.6 _(2.7)
OPT-125M TINT	32	82.3 _(2.7)	69.3 _(0.9)	73.7 _(0.8)	75.7 _(1.9)	72.3 _(1.2)	83.2 _(1.0)	78.2 _(0.2)	76.4 _(0.7)
<i>With Calibration</i>									
OPT-125M	0	64.0	66.0	53.0	54.5	52.5	55.5	58.0	57.6
OPT-1.3B	0	73.5	61.5	57.5	53.0	54.5	79.5	61.0	62.9
OPT-125M DYNA. EVAL	0	62.5	66.0	60.5	53.5	54.0	56.5	74.5	61.1
OPT-125M TINT	0	64.0	66.0	56.5	59.0	53.5	62.0	66.5	61.1
OPT-125M	32	83.5 _(2.4)	40.7 _(10.4)	50.8 _(0.8)	67.7 _(4.1)	57.7 _(10.8)	79.2 _(8.4)	56.0 _(8.1)	62.2 _(2.7)
OPT-1.3B	32	51.8 _(1.9)	66.2 _(3.1)	93.7 _(1.0)	82.8 _(2.8)	91.3 _(1.9)	83.5 _(2.5)	92.0 _(2.9)	80.2 _(0.7)
OPT-125M DYNA. EVAL	32	87.2 _(0.2)	67.2 _(0.6)	72.8 _(5.9)	73.3 _(2.6)	66.7 _(7.4)	81.5 _(3.7)	70.3 _(2.1)	74.1 _(2.9)
OPT-125M TINT	32	85.3 _(1.9)	67.3 _(0.6)	71.8 _(3.8)	70.7 _(1.9)	63.7 _(0.2)	83.5 _(1.6)	77.5 _(1.2)	74.3 _(1.4)

TINT remains comparable to dynamic evaluation. Additionally, we find that TINT outperforms or is on par with a similarly sized pre-trained model (OPT-1.3B) except in the calibrated few-shot setting. This suggests that the capabilities of existing pre-trained models may be understood via the simulation of smaller auxiliary models. Please refer to Appendix L for more experiment details.

4 Related Work

Fast Weight Programmers (FWPs) enable input-dependent weight updates during inference. Ba et al. [3] connect self-attention and FWPs, and follow-up works [37, 19] show the efficacy of self-attention layers to update linear and recurrent networks during inference. Clark et al. [10] added Fast Weights Layers (FWL) to a frozen pre-trained model, and efficiently fine-tune FWL as the model processes the sequence.

Alternative Explanations for ICL: Some works study ICL using a Bayesian framework. Xie et al. [50] model pretraining data as a mixture of HMMs and cast ICL identifying one such component. Hahn and Goyal [15] later modeled language as a compositional grammar, and propose ICL as a composition of operations. On the other hand, careful experiments in Chan et al. [7] show that data distributional properties (e.g. Zipf’s law) drive in-context learning in transformers.

5 Discussion

We present a parameter-efficient construction TINT capable of simulating gradient descent on an internal transformer model during inference. Using fewer than 2 billion parameters, it can simulate fine-tuning a 125 million transformer (e.g., GPT-2) internally, dramatically reducing the scale required by previous works. Language modeling and in-context learning experiments demonstrate that the efficient approximations still allow the TINT to fine-tune the model. Our work emphasizes that the inference behavior of complex models may rely on the training dynamics of smaller models. As such, the existence of TINT has strong implications for interpretability and AI alignment research.

While our work represents a significant improvement over previous simulations in terms of auxiliary model complexity, similar to prior research in this area, our insights into existing pre-trained models are limited. Furthermore, we have not yet examined potential biases that may arise in the auxiliary models due to one-step gradient descent. We plan to investigate these aspects in future work.

References

- [1] Ekin Akyurek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- [2] Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization in large language models. *arXiv preprint arXiv:2207.04901*, 2022.
- [3] Jimmy Ba, Geoffrey Hinton, Volodymyr Mnih, Joel Z. Leibo, and Catalin Ionescu. Using fast weights to attend to the recent past, 2016.
- [4] Satwik Bhattamishra, Kabir Ahuja, and Navin Goyal. On the ability and limitations of transformers to recognize formal languages. *arXiv preprint arXiv:2009.11264*, 2020.
- [5] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shrivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR, 2023.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [7] Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond, James McClelland, and Felix Hill. Data distributional properties drive emergent in-context learning in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891, 2022.
- [8] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [9] Bilal Chughtai, Lawrence Chan, and Neel Nanda. A toy model of universality: Reverse engineering how networks learn group operations. *arXiv preprint arXiv:2302.03025*, 2023.
- [10] Kevin Clark, Kelvin Guu, Ming-Wei Chang, Panupong Pasupat, Geoffrey Hinton, and Mohammad Norouzi. Meta-learning fast weight language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9751–9757, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.emnlp-main.661>.
- [11] Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Zhifang Sui, and Furu Wei. Why can gpt learn in-context? language models secretly perform gradient descent as meta-optimizers, 2022.
- [12] Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable creation in self-attention mechanisms. In *International Conference on Machine Learning*, pages 5793–5831. PMLR, 2022.
- [13] N Elhage, N Nanda, C Olsson, T Henighan, N Joseph, B Mann, A Askell, Y Bai, A Chen, T Conerly, et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021.
- [14] Angeliki Giannou, Shashank Rajput, Jy yong Sohn, Kangwook Lee, Jason D. Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers, 2023.
- [15] Michael Hahn and Navin Goyal. A theory of emergent in-context learning as implicit structure induction. *arXiv preprint arXiv:2303.07971*, 2023.
- [16] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.

- [17] Ari Holtzman, Peter West, Vered Shwartz, Yejin Choi, and Luke Zettlemoyer. Surface form competition: Why the highest probability answer isn’t always right. *arXiv preprint arXiv:2104.08315*, 2021.
- [18] Minqing Hu and Bing Liu. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 168–177, 2004.
- [19] Kazuki Irie, Imanol Schlag, Róbert Csordás, and Jürgen Schmidhuber. Going beyond linear transformers with recurrent fast weight programmers. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=ot20RiBqTa1>.
- [20] Ben Krause, Emmanuel Kahembwe, Iain Murray, and Steve Renals. Dynamic evaluation of transformer language models. *arXiv preprint arXiv:1904.08378*, 2019.
- [21] Ananya Kumar, Ruqi Shen, Sébastien Bubeck, and Suriya Gunasekar. How to fine-tune vision models with sgd, 2022.
- [22] David Lindner, János Kramár, Matthew Rahtz, Thomas McGrath, and Vladimir Mikulik. Tracr: Compiled transformers as a laboratory for interpretability. *arXiv preprint arXiv:2301.05062*, 2023.
- [23] Bingbin Liu, Jordan T. Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers learn shortcuts to automata. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=De4FYqjFueZ>.
- [24] Peter J Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. Generating wikipedia by summarizing long sequences. *arXiv preprint arXiv:1801.10198*, 2018.
- [25] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*, 2016.
- [26] Neel Nanda, Lawrence Chan, Tom Liberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023.
- [27] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- [28] Bo Pang and Lillian Lee. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04)*, pages 271–278, 2004.
- [29] Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05)*, pages 115–124, 2005.
- [30] Jorge Perez, Pablo Barcelo, and Javier Marinkovic. Attention is turing-complete. *Journal of Machine Learning Research*, 22(75):1–35, 2021. URL <http://jmlr.org/papers/v22/20-302.html>.
- [31] Ofir Press, Noah A Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables input length extrapolation. *arXiv preprint arXiv:2108.12409*, 2021.
- [32] Jorge Pérez, Javier Marinković, and Pablo Barceló. On the turing completeness of modern neural network architectures. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HyGBdoOqFm>.
- [33] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

- [34] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [35] Nikunj Saunshi, Sadhika Malladi, and Sanjeev Arora. A mathematical exploration of why language models help solve downstream tasks. *arXiv preprint arXiv:2010.03648*, 2020.
- [36] Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- [37] Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight memory systems. *CoRR*, abs/2102.11174, 2021. URL <https://arxiv.org/abs/2102.11174>.
- [38] Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- [39] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642, 2013.
- [40] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021.
- [41] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [43] Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent, 2022.
- [44] Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *NeurIPS ML Safety Workshop*, 2022. URL <https://openreview.net/forum?id=rvi3Wa768B->.
- [45] Colin Wei, Yining Chen, and Tengyu Ma. Statistically meaningful approximation: a case study on approximating turing machines with transformers. *CoRR*, abs/2107.13163, 2021. URL <https://arxiv.org/abs/2107.13163>.
- [46] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=yzkSU5zdwD>. Survey Certification.
- [47] Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In *International Conference on Machine Learning*, pages 11080–11090. PMLR, 2021.
- [48] Janyce Wiebe, Theresa Wilson, and Claire Cardie. Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 39:165–210, 2005.
- [49] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018.
- [50] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=RdJVFCHjUMI>.

- [51] Shunyu Yao, Binghui Peng, Christos Papadimitriou, and Karthik Narasimhan. Self-attention networks can process bounded hierarchical languages. *arXiv preprint arXiv:2105.11115*, 2021.
- [52] Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [53] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- [54] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28, 2015.

Contents

1	Introduction	1
2	Our Construction	2
2.1	Overview	2
2.2	Key Components	3
2.3	Operating on an auxiliary model with prefix embeddings	4
2.4	Stacking in prefix-tokens, H_{sim} -split linear operations and Linear attention	4
3	Verification of TINT	5
4	Related Work	6
5	Discussion	6
A	Overview of remaining modifications	12
A.1	First order gradients for layer normalization	12
A.2	Backpropagation through Attention Value Vectors	12
A.3	Parameter sharing in the TINT	13
B	Additional Related works	13
C	Additional Notations	14
C.1	Simulating Multiplication from [1]	16
D	Linear layer	16
D.1	H_{sim} -split operation	18
E	Self-attention layer	20
E.1	Approximate auxiliary self-attention backpropagation	24
E.2	Proofs of theorems and gradient definitions	25
F	Layer normalization	28
F.1	Additional definitions	30
F.2	Proof of theorems and gradient definitions	30
G	Activation layer	32
G.1	Proofs of theorems	34
H	Language model head	35
I	Parameter sharing	36
J	Additional modules	36

J.1	Root mean square normalization (RMSnorm)	36
J.2	Attention variants	37
J.3	Gated linear units (GLUs)	37
K	Construction of other variants of pre-trained models	40
L	Experiments	40
L.1	Experimental Setup	40
L.2	Additional Experiments	41
M	Broader Impacts	42

A Overview of remaining modifications

A.1 First order gradients for layer normalization

Below, we show that computing exact gradients for layer normalization is expensive, so we efficiently approximate backpropagation by computing the dominating term.

Definition A.1. [Layer Normalization] Define a normalization function $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that performs $f(\mathbf{x}) = (\mathbf{x} - \mu)/\sigma$, where μ and σ are the mean and standard deviation of \mathbf{x} , respectively. Then, layer normalization with parameters $\gamma, \mathbf{b} \in \mathbb{R}^{D_{\text{aux}}}$ takes as input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ and outputs $\mathbf{y} \in \mathbb{R}^{D_{\text{aux}}}$, which is computed as $\mathbf{z} = f(\mathbf{x}), \mathbf{y} = \gamma \odot \mathbf{z} + \mathbf{b}$.

Definition A.2. [Exact Gradient for Layer Normalization] Using notations in Definition A.1, given the gradient of the loss w.r.t the output of the Layer Normalization $\partial_{\mathbf{y}}$, backpropagation computes $\partial_{\mathbf{x}}$ as

$$\partial_{\mathbf{x}} = (\partial_{\mathbf{z}} - D_{\text{aux}}^{-1} \sum_{i=1}^{D_{\text{aux}}} \partial_{z_i} - \langle \partial_{\mathbf{z}}, \mathbf{z} \rangle \mathbf{z}) / \sigma \quad \partial_{\mathbf{z}} = \gamma \odot \partial_{\mathbf{y}}.$$

Exact backpropagation is expensive because $\langle \partial_{\mathbf{z}}, \mathbf{z} \rangle \mathbf{z}$ requires using at least two sequential MLPs. We thus approximate it with a first-order Taylor expansion, which is entry-wise close to the true gradient.

Definition A.3. [ϵ -approximate Layer Normalization Gradient] With notations defined above, this layer takes $\partial_{\mathbf{y}}, \mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ as input and outputs $\widehat{\partial_{\mathbf{x}}} = \frac{1}{\epsilon} (f(\mathbf{x} + \epsilon \gamma \odot \partial_{\mathbf{y}}) - f(\mathbf{x}))$.

Theorem A.4 (Informal, c.f. Thm F.1). *With bounded ℓ_2 -norms of $\mathbf{x}, \partial_{\mathbf{y}}, \gamma, \mathbf{b}$, $\|\partial_{\mathbf{x}} - \widehat{\partial_{\mathbf{x}}}\|_{\infty} \leq \mathcal{O}(\epsilon)$.*

This approximation only works for symmetric Jacobian, so it does not apply when backpropagating the linear and self-attention layers. We use a TINT module with 2 linear layers, separated by Group Normalization [49], to compute $\widehat{\partial_{\mathbf{x}}}$, resulting in a $4 \times$ parameter reduction.

A.2 Backpropagation through Attention Value Vectors

For simplicity, we show a single head self-attention layer (multi-head attention is in Appendix E). The self-attention in the auxiliary model is the same as in TINT (Definition 2.3) without a position vector.

Definition A.5 (Auxiliary model softmax self-attention). A self-attention layer with parameters $\{\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V\}$ takes a sequence $\{\mathbf{x}_t\}_{t \leq T_{\text{aux}}}$ and outputs a sequence $\{\mathbf{y}_t\}_{t \leq T_{\text{aux}}}$, such that

$$\mathbf{y}_t = \sum_j a_{t,j} \mathbf{v}_j, \quad \text{with } a_{t,j} = \text{softmax}(\mathbf{K} \mathbf{q}_t)_j, \quad \mathbf{q}_t = \mathbf{W}_Q \mathbf{x}_t, \quad \mathbf{k}_t = \mathbf{W}_K \mathbf{x}_t, \quad \mathbf{v}_t = \mathbf{W}_V \mathbf{x}_t,$$

for all $t \leq T_{\text{aux}}$, and $\mathbf{K} \in \mathbb{R}^{T_{\text{aux}} \times D_{\text{aux}}}$ defined with rows $\{\mathbf{k}_t\}_{t=1}^{T_{\text{aux}}}$.

Table 2: Language modeling results on WIKITEXT-103. We use 30%, 50%, 70% and 90% of sequences for training in dynamic eval and TINT and the rest of the sequence for evaluation. TINT improves upon the auxiliary model perplexities by 0.3 – 0.7 absolute on average. The small perplexity difference between the TINT and dynamic evaluation suggests that the approximations introduced in the descent algorithm (Section 2.2) can still effectively fine-tune the auxiliary model.

Training proportion	GPT2				OPT-125m			
	30%	50%	70%	90%	30%	50%	70%	90%
VANILLA MODEL	25.6	24.9	24.5	23.3	29.6	28.8	28.0	28.0
DYNA. EVAL	24.9	24.0	23.5	22.2	29.0	28.2	27.4	27.4
TINT	25.1	24.3	23.8	22.6	29.3	28.4	27.5	27.4

Definition A.6 (Exact gradient for softmax self-attention). Given the gradients of the loss w.r.t the output sequence $\{\partial_{\mathbf{y}_t}\}_{t=1}^T$, backpropagation computes $\{\partial_{\mathbf{x}_t}\}_{t=1}^T$, with

$$\begin{aligned}\partial_{\mathbf{x}_t} &= \mathbf{W}_Q^\top \partial_{\mathbf{q}_t} + \mathbf{W}_K^\top \partial_{\mathbf{k}_t} + \mathbf{W}_V^\top \partial_{\mathbf{v}_t}, & \partial_{\mathbf{v}_t} &= \sum_j a_{j,t} \partial_{\mathbf{y}_j}, \\ \partial_{\mathbf{q}_t} &:= \sum_j a_{t,j} ((\partial_{\mathbf{y}_t})^\top \mathbf{v}_j) [\mathbf{k}_j - \sum_{j'} a_{t,j'} \mathbf{k}_{j'}], & \partial_{\mathbf{k}_t} &:= \sum_j a_{t,j} ((\partial_{\mathbf{y}_t})^\top (\mathbf{v}_j - \sum_{j'} a_{t,j'} \mathbf{v}_{j'})) \mathbf{q}_j\end{aligned}$$

for all $t \leq T_{\text{aux}}$, with $\mathbf{K} \in \mathbb{R}^{T_{\text{aux}} \times D_{\text{aux}}}$ defined with rows $\{\mathbf{k}_t\}_{t=1}^{T_{\text{aux}}}$.

The computation of $\partial_{\mathbf{q}_t} := \sum_j a_{t,j} ((\partial_{\mathbf{y}_t})^\top \mathbf{v}_j) [\mathbf{k}_j - \sum_{j'} a_{t,j'} \mathbf{k}_{j'}]$ (and similarly $\partial_{\mathbf{k}_t}$) requires at least 2 self-attention layers and an MLP layer, as we must compute and multiply attention scores $a_{t,j}$ and $(\partial_{\mathbf{y}_t})^\top \mathbf{v}_j$ before computing $\partial_{\mathbf{q}_t}$. Thus, we only update the self-attention using the gradients w.r.t. \mathbf{v}_t .

Definition A.7 (Approximate Self-Attention Backpropagation). With the notations defined above, this layer takes a sequence $\{\partial_{\mathbf{y}_t} \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and $\{\mathbf{x}_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ as input and outputs $\{\widehat{\partial_{\mathbf{x}_t}}\}_{t \leq T}$, with $\widehat{\partial_{\mathbf{x}_t}} = \mathbf{W}_V^\top \partial_{\mathbf{v}_t}$, where $\partial_{\mathbf{v}_t} = \sum_j a_{j,t} \partial_{\mathbf{y}_j}$.

We formally show that when the attention head for each position pays a lot of attention to a single token (i.e., behaves like hard attention [30]), $\widehat{\partial_{\mathbf{x}_t}}$ is entry-wise close to $\partial_{\mathbf{x}_t}$ for all t . Computing $\{\widehat{\partial_{\mathbf{x}_t}}\}_{t=1}^T$ instead of $\{\partial_{\mathbf{x}_t}\}_{t=1}^T$ induces a $5 \times$ parameter reduction.

Theorem A.8 (Informal, c.f. Theorem E.4). *If on input sequence $\{\mathbf{x}_t\}_{t \leq T_{\text{aux}}}$, the attention scores are ε -close to a hard-attention at each position, then for all t , $\|\partial_{\mathbf{x}_t} - \widehat{\partial_{\mathbf{x}_t}}\| \leq \mathcal{O}(\varepsilon)$.*

A.3 Parameter sharing in the TINT

Consider the self-attention layer (Appendix A.2). The relevant TINT module performs linear operations with \mathbf{W}_Q , \mathbf{W}_K , \mathbf{W}_V to compute query, key, and value vectors at each position t (Definition 2.2) and hence can be simulated with the Linear Forward module (Section 2.4). We additionally leverage parameter sharing to apply a single Linear Forward module for each of the three computations, changing only the prefix embeddings to correspond to \mathbf{W}_Q , \mathbf{W}_K , or \mathbf{W}_V . Applying the same structure to feed-forward linear layers results in a $4 \times$ reduction in the number of necessary modules (Appendix I).

B Additional Related works

Interpretability: Mechanistic interpretability works reverse-engineer the algorithms simulated by these models [13, 27, 44, 26, 9]. These works study local patterns, e.g. activations and attention heads, to derive interpretable insights. Other works [47, 22] use declarative programs to algorithmically describe transformer models.

Transformer Expressivity: Perez et al. [30], Pérez et al. [32] show that Transformers with hard attention are Turing complete, with Wei et al. [45] showing statistically meaningful transformer constructions for Turing machines for statistical learnability. In Section 2.3, we point out that this scheme often results in gigantic constructions. To understand the behavior of moderate sized models,

other works have investigated specific classes of algorithms, e.g. bounded-depth Dyck languages [51], modular prefix sums [2], adders [26], regular languages [4], and sparse logical predicates [12]. Liu et al. [23] provide a unified theory to understand automata-like mechanisms within transformers.

C Additional Notations

We differentiate the parameters of the auxiliary model and TINT by using an explicit superscript TINT for TINT parameters, for example, the weights of a linear layer in TINT will be represented by \mathbf{W}^{TINT} . We use two operations throughout: SPLIT_h and VECTORIZE . Function $\text{SPLIT}_h : \mathbb{R}^d \rightarrow \mathbb{R}^{h \times \lfloor d/h \rfloor}$ takes an input $\mathbf{x} \in \mathbb{R}^d$ and outputs H equal splits of \mathbf{x} , for any arbitrary dimension d . Function $\text{VECTORIZE} : \mathbb{R}^{h \times d} \rightarrow \mathbb{R}^{dh}$ concatenates the elements of a sequence $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i \leq h}$ into one single vector, for any arbitrary d and h . Recall that for a matrix \mathbf{A} , \mathbf{a}_j refers to its j th row, and for any vector \mathbf{b} , b_j refers to its j th element. However, at a few places in the appendix, for typographical reasons, for a matrix \mathbf{A} , we have also used $(\mathbf{A})_j$ to refer to its j th row, and for any vector \mathbf{b} , $(\mathbf{b})_j$ to refer to its j th element.

TINTAttention Module We modify the usual attention module to include the position embeddings $\{\mathbf{p}_i^{\text{TINT}} \in \mathbb{R}^{T_{\text{sim}}}\}_{i \leq T_{\text{sim}}}$. In usual self-attention modules, the query, key, and value vectors at each position are computed by token-wise linear transformations of the input embeddings. In TINT’s Attention Module, we perform additional linear transformations on the position embeddings, using parameters $\mathbf{W}_Q^p, \mathbf{W}_K^p, \mathbf{W}_V^p$, and decision vectors $\lambda^Q, \lambda^K, \lambda^V \in \mathbb{R}^{H_{\text{sim}}}$ decide whether to add these transformed position vectors to the query, key, and value vectors of different attention heads. The following definition generalizes single-head attention (Definition 2.3). For the following definition, we use $\hat{\mathbf{e}}$ to represent input sequence and $\tilde{\mathbf{e}}$ to represent the output sequence: we introduce these general notations below to avoid confusion with the notations for token and prefix embeddings for TINT illustrated in Figure 1.

Definition C.1 (TINT’s self-attention with H_{sim} heads). For parameters $\{\mathbf{W}_Q^{\text{TINT}}, \mathbf{W}_K^{\text{TINT}}, \mathbf{W}_V^{\text{TINT}} \in \mathbb{R}^{D_{\text{sim}} \times D_{\text{sim}}}\}$, $\{\mathbf{b}_Q^{\text{TINT}}, \mathbf{b}_K^{\text{TINT}}, \mathbf{b}_V^{\text{TINT}} \in \mathbb{R}^{D_{\text{sim}}}\}$, $\{\mathbf{W}_Q^p, \mathbf{W}_K^p, \mathbf{W}_V^p \in \mathbb{R}^{T_{\text{sim}} \times D_{\text{sim}}/H_{\text{sim}}}\}$ and $\{\lambda^Q, \lambda^K, \lambda^V \in \mathbb{R}^{H_{\text{sim}}}\}$, TINT self-attention with H_{sim} attention heads and a function $f_{\text{attn}} : \mathbb{R}^{T_{\text{sim}}} \rightarrow \mathbb{R}^{T_{\text{sim}}}$ takes a sequence $\{\hat{\mathbf{e}}_t \in \mathbb{R}^{D_{\text{sim}}}\}_{t \leq T_{\text{sim}}}$ as input and outputs $\{\tilde{\mathbf{e}}_t \in \mathbb{R}^{D_{\text{sim}}}\}_{t \leq T_{\text{sim}}}$, with

$$\begin{aligned} \tilde{\mathbf{e}}_t &= \text{VECTORIZE}(\{\sum_{j \leq T_{\text{sim}}} a_{t,j}^h \tilde{\mathbf{v}}_j^h\}_{h \leq H_{\text{sim}}}), \text{ with } a_{t,j}^h = f_{\text{attn}}(\tilde{\mathbf{K}}^h \tilde{\mathbf{q}}_t^h)_j \\ \tilde{\mathbf{q}}_t^h &= \text{SPLIT}_H(\mathbf{q}_t)_h + \lambda_h^Q \mathbf{W}_Q^p \mathbf{p}_t^{\text{TINT}}; \quad \tilde{\mathbf{k}}_t^h = \text{SPLIT}_H(\mathbf{k}_t)_h + \lambda_h^K \mathbf{W}_K^p \mathbf{p}_t^{\text{TINT}}; \\ \tilde{\mathbf{v}}_t^h &= \text{SPLIT}_H(\mathbf{v}_t)_h + \lambda_h^V \mathbf{W}_V^p \mathbf{p}_t^{\text{TINT}}. \end{aligned}$$

Here, $\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t$ denote the query, key, and value vectors at each position t , computed as $\mathbf{W}_Q^{\text{TINT}} \hat{\mathbf{e}}_t + \mathbf{b}_Q^{\text{TINT}}, \mathbf{W}_K^{\text{TINT}} \hat{\mathbf{e}}_t + \mathbf{b}_K^{\text{TINT}}$, and $\mathbf{W}_V^{\text{TINT}} \hat{\mathbf{e}}_t + \mathbf{b}_V^{\text{TINT}}$ respectively. $\tilde{\mathbf{K}}^h \in \mathbb{R}^{T_{\text{sim}} \times D_{\text{sim}}/H_{\text{sim}}}$ is defined with its rows as $\{\tilde{\mathbf{k}}_t^h\}_{t \leq T_{\text{sim}}}$ for all $h \leq H_{\text{sim}}$.

f_{attn} can be either linear or softmax function.

Bounded parameters and input sequence: We define a linear self-attention layer to be B_w -bounded, if the ℓ_2 norms of all the parameters are bounded by B_w . Going by Definition C.1, this implies

$$\begin{aligned} \max\{\|\mathbf{W}_Q^{\text{TINT}}\|_2, \|\mathbf{W}_K^{\text{TINT}}\|_2, \|\mathbf{W}_V^{\text{TINT}}\|_2\} &\leq B_w, \quad \max\{\|\mathbf{b}_Q^{\text{TINT}}\|_2, \|\mathbf{b}_K^{\text{TINT}}\|_2, \|\mathbf{b}_V^{\text{TINT}}\|_2\} \leq B_w \\ \max\{\|\mathbf{W}_Q^p\|_2, \|\mathbf{W}_K^p\|_2, \|\mathbf{W}_V^p\|_2\} &\leq B_w, \quad \max\{\|\lambda^Q\|_2, \|\lambda^K\|_2, \|\lambda^V\|_2\} \leq B_w. \end{aligned}$$

Furthermore, we define an input sequence $\{\hat{\mathbf{e}}_t\}_{t \leq T_{\text{sim}}}$ to be B_x -bounded, if $\|\hat{\mathbf{e}}_t\|_2 \leq B_x$ for all t .

Recall from the main paper (Section 2.4), we used Linear TINT Self-Attention layer to represent the linear operations of the auxiliary model. In the following theorem, we show that a linear attention layer can be represented as a softmax attention layer that uses an additional attention head and an extra token \mathbf{u} , followed by a linear layer. Therefore, replacing softmax attention with linear attention does not deviate too far from the canonical transformer. We use the Linear TINT Self-Attention layers in several places throughout the model.

Theorem C.2. For any $B_w > 0$, consider a B_w -bounded linear self-attention layer that returns $\{\tilde{\mathbf{e}}_t^{linear} \in \mathbb{R}_{sim}^D\}_{t \leq T_{sim}}$ on any input $\{\hat{\mathbf{e}}_t \in \mathbb{R}_{sim}^D\}_{t \leq T_{sim}}$. Consider a softmax self-attention layer with $2H_{sim}$ attention heads and an additional token $\mathbf{u} \in \mathbb{R}^{2D_{sim}}$ such that for any B_x -bounded input $\{\hat{\mathbf{e}}_t\}_{t \leq T_{sim}}$, it takes a modified input sequence $\{\bar{\mathbf{e}}_1, \dots, \bar{\mathbf{e}}_{T_{sim}}, \mathbf{u}\}$, and returns $\{\tilde{\mathbf{e}}_t^{softmax} \in \mathbb{R}^{2D_{sim}}\}_{t \leq T_{sim}}$. Each modified input token $\bar{\mathbf{e}}_t \in \mathbb{R}^{2D_{sim}}$ is obtained by concatenating additional 0s to $\hat{\mathbf{e}}_t$. Then, for any $B_x > 0$, and $\epsilon \leq \mathcal{O}(T_{sim}^{-2} B_w^{-5} B_x^{-5})$, there exists $\mathbf{W}_O \in \mathbb{R}^{D_{sim} \times 2D_{sim}}$ and such a softmax self-attention layer such that

$$\left\| \mathbf{W}_O \tilde{\mathbf{e}}_t^{softmax} - \tilde{\mathbf{e}}_t^{linear} \right\|_2 \leq \mathcal{O}(\sqrt{\epsilon}),$$

for all $t \leq T_{sim}$.

Proof. Consider an input sequence $\{\mathbf{x}_t\}_{t \leq T_{sim}}$. Let the attention scores of any linear head $h \leq H_{sim}$ in the linear attention layer be given by $\{a_{t,j}^h\}_{j \leq T_{sim}}$, at any given position t . Additionally, let the value vectors for the linear attention be given by \mathbf{v}_t . To repeat our self-attention definition, the output of the attention layer at any position t is given by $\text{VECTORIZE}(\{\tilde{\mathbf{e}}_t^{linear,h}\}_{h \leq H_{sim}})$, where

$$\tilde{\mathbf{e}}_t^{linear,h} = \sum_{j \leq T_{sim}} a_{t,j}^h \mathbf{v}_j^h.$$

Under our assumption, B_w denotes the maximum ℓ_2 norm of all the parameters in the linear self-attention layer and B_x the maximum ℓ_2 norm in the input sequence, i.e. $\max_{t \leq T_{sim}} \|\mathbf{x}_t\|_2 \leq B_x$. With a simple application of Cauchy-Schwartz inequality, we can show that $\max_{j \leq T_{sim}} |a_{t,j}^h| \leq \mathcal{O}(B_w^2 B_x^2)$, and $\max_{t \leq T_{sim}} \|\mathbf{v}_t^h\|_2 \leq \mathcal{O}(B_w B_x)$.

For $\epsilon \leq \mathcal{O}(T_{sim}^{-10/9} B_w^{-40/9} B_x^{-40/9})$, we can then use Lemma C.3 to represent for each $t, j \leq T_{sim}$,

$$\begin{aligned} a_{t,j}^h &= \frac{\epsilon^{-3} e^{\epsilon a_{t,j}^h}}{\sum_{t' \leq T_{sim}} e^{\epsilon a_{t,t'}^h} + e^{-2 \log \epsilon}} - \epsilon^{-1} + \mathcal{O}(\epsilon(T_{sim} + a_{t,j}^h)) \\ &:= \epsilon^{-3} \text{softmax}(\{\epsilon a_{t,1}^h, \epsilon a_{t,2}^h, \dots, \epsilon a_{t,T_{sim}}^h, -2 \log \epsilon\})_j - \epsilon^{-1} + \mathcal{O}(\epsilon^{0.9}). \end{aligned}$$

Softmax attention construction: We define \mathbf{u} , and the query and key parameters of the softmax attention layer such that for the first H_{sim} attention heads, the query-key dot products for all the attention heads between any pairs $\{(\bar{\mathbf{e}}_t, \bar{\mathbf{e}}_j)\}_{t,j \leq T_{sim}}$ is given by $\{\epsilon a_{t,j}^h\}_{h \leq H_{sim}}$, while being $-2 \log \epsilon$ between \mathbf{u} and any token $\bar{\mathbf{e}}_t$, with $t \leq T_{sim}$. For the rest of H_{sim} attention heads, the attention scores are uniformly distributed across all pairs of tokens (attention score between any pair of tokens is given by $\frac{1}{T_{sim}+1}$).

We set the value parameters of the softmax attention layer such that at any position $t \leq T_{sim}$, the value vector is given by $\text{VECTORIZE}(\{\epsilon^{-3} \mathbf{v}_t, \mathbf{v}_t\})$. The value vector returned for \mathbf{u} contains all 0s.

Softmax attention computation: Consider an attention head $h \leq H_{sim}$ in the softmax attention layer now. The output of the attention head at any position $t \leq T_{sim}$ is given by

$$\begin{aligned} \tilde{\mathbf{e}}_t^{softmax,h} &= \sum_{j \leq T_{sim}} \text{softmax}(\{\epsilon a_{t,1}^h, \epsilon a_{t,2}^h, \dots, \epsilon a_{t,T_{sim}}^h, -2 \log \epsilon\})_j \epsilon^{-3} \mathbf{v}_j^h \\ &= \sum_{j \leq T_{sim}} (a_{t,j}^h + \epsilon^{-1} + \mathcal{O}(\epsilon^{0.9})) \mathbf{v}_j^h. \end{aligned}$$

This has an additional $\sum_{j \leq T_{sim}} (\epsilon^{-1} + \mathcal{O}(\epsilon^{0.9})) \mathbf{v}_j^h$, compared to $\tilde{\mathbf{e}}_t^{linear,h}$. However, consider the output of the attention head $H_{sim} + h$ at the same position:

$$\tilde{\mathbf{e}}_t^{softmax, H_{sim}+h} = \frac{1}{T_{sim}+1} \sum_{j \leq T_{sim}} \mathbf{v}_j^h.$$

Hence, we can use the output matrix \mathbf{W}_O to get $\tilde{\mathbf{e}}_t^{softmax,h} - \frac{T_{sim}+1}{\epsilon} \tilde{\mathbf{e}}_t^{softmax, H_{sim}+h} = \sum_{j \leq T_{sim}} (a_{t,j}^h + \mathcal{O}(\epsilon^{0.9})) \mathbf{v}_j^h$. The additional term $\mathcal{O}(\epsilon^{0.9}) \sum_{j \leq T_{sim}} \mathbf{v}_j^h$ can be further shown to be $\mathcal{O}(\epsilon^{0.5})$ small with the assumed bound of ϵ , since each \mathbf{v}_j^h is atmost $\mathcal{O}(B_w B_x)$ in ℓ_2 norm with a Cauchy Schwartz inequality. \square

Lemma C.3. For $\epsilon > 0$, $B > 0$, and a sequence $\{a_1, a_2, \dots, a_T\}$ with each $a_i \in \mathbb{R}$ and $|a_i| \leq B$, the following holds true for all $i \leq T$,

$$\frac{\epsilon^{-3}e^{\epsilon a_i}}{\sum_{t' \leq T} e^{\epsilon a_{t'}} + e^{-2 \log \epsilon}} = a_i + \frac{1}{\epsilon} + \mathcal{O}(\epsilon^{0.9}),$$

provided $\epsilon \leq \mathcal{O}(T^{-10/9} B^{-20/9})$.

Proof. We will use the following first-order Taylor expansions:

$$e^x = 1 + x + \mathcal{O}(x^2). \quad (1)$$

$$\frac{1}{1+x} = 1 - \mathcal{O}(x). \quad (2)$$

Hence, for any $x \ll 1$, $x \approx e^x - 1$.

Simplifying the L.H.S. of the desired bound, we have

$$\frac{\epsilon^{-3}e^{\epsilon a_i}}{\sum_{t' \leq T} e^{\epsilon a_{t'}} + e^{-2 \log \epsilon}} = \frac{\epsilon^{-3}(1 + \epsilon a_i + \mathcal{O}(\epsilon^2 a_i^2))}{\sum_{t' \leq T} (1 + \epsilon a_{t'} + \mathcal{O}(\epsilon^2 a_{t'}^2)) + e^{-2 \log \epsilon}} \quad (3)$$

$$= \frac{\epsilon^{-1} + a_i + \mathcal{O}(\epsilon a_i^2)}{\sum_{t' \leq T} (\epsilon^2 + \epsilon^3 a_{t'} + \mathcal{O}(\epsilon^4 a_{t'}^2)) + 1} \quad (4)$$

$$= (\epsilon^{-1} + a_i + \mathcal{O}(\epsilon a_i^2)) (1 + \mathcal{O}(\epsilon^2 T)) \quad (5)$$

$$= \epsilon^{-1} + a_i + \mathcal{O}(\epsilon T + a_i^2 T \epsilon^2 + a_i^2 T \epsilon^3 + \epsilon a_i^2) = \epsilon^{-1} + a_i + \mathcal{O}(\epsilon^{0.9}).$$

We used Taylor expansion of exponential function (Equation (1)) in Equation (3) to get Equation (4), and Taylor expansion of inverse function (Equation (2)) to get Equation (5) from Equation (4). Furthermore, with the lower bound assumption on ϵ , $\sum_{t' \leq T} (\epsilon^2 + \epsilon^3 a_{t'} + \mathcal{O}(\epsilon^4 a_{t'}^2))$ can be shown to be at most $3\epsilon^2 T$, which amounts to $\mathcal{O}(\epsilon^2 T)$ error in Equation (5). The final error bound has again been simplified using the lower bound assumption on ϵ . \square

C.1 Simulating Multiplication from [1]

We refer to the multiplication strategy of [1] at various places.

Lemma C.4. [Lemma 4 in [1]] The *GeLU* [16] nonlinearity can be used to perform multiplication: specifically,

$$\sqrt{\pi/2}(\text{GeLU}(x+y) - \text{GeLU}(y)) = xy + \mathcal{O}(x^3 y^3).$$

Thus, to represent an element-wise product or a dot product between two sub-vectors in a token embedding, we can use a MLP with a *GeLU* activation.

D Linear layer

In the main paper, we defined the linear layer without the bias term for simplicity (Definition 2.2). In this section, we will redefine the linear layer with the bias term and present a comprehensive construction of the Linear Forward module.

Definition D.1 (Linear layer). For a weight $\mathbf{W} \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$ and bias $\mathbf{b} \in \mathbb{R}^{D_{\text{aux}}}$, a linear layer takes $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ as input and outputs $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$.

In the discussions below, we consider a linear layer in the auxiliary model with parameters $\{\mathbf{W}, \mathbf{b}\}$ that takes in input sequence $\mathbf{x}_1, \dots, \mathbf{x}_{T_{\text{aux}}}$ and outputs $\mathbf{y}_1, \dots, \mathbf{y}_{T_{\text{aux}}}$, with $\mathbf{y}_t = \mathbf{W}\mathbf{x}_t + \mathbf{b}$ for each $t \leq T_{\text{aux}}$. Since this involves a token-wise operation, we will present our constructed modules with a general token position t and the prefix tokens $\{\mathbf{v}_j\}$.

TINT Linear Forward module Continuing our discussion from Section 2.4, we represent S stacked rows of \mathbf{W} as a prefix embedding. In addition, we store the bias \mathbf{b} in the first prefix embedding (\mathbf{v}_1).

Using a set of S' unique attention heads in a TINT attention module (Definition C.1), we copy the bias \mathbf{b} to respective token embeddings and use a TINT linear layer to add the biases to the final output.

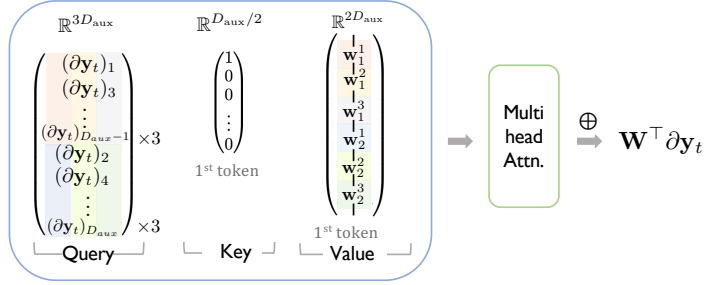


Figure 3: TINT simulates the backward pass of a linear layer as a H -head attention layer ($H = 6$ pictured), with the gradient of the loss w.r.t. linear layer output (∂y_t) as the query, the positional one-hot vector of prefix embeddings as the key, and the parameters of the auxiliary model stored in the prefix embeddings as the value. Similar to the Linear Forward module (Figure 2), we distribute the dot product computations across all attention heads by sharding the vectors into S' ($S' = 3$ here) parts. We omitted the identical transformation for query, and value matrices, and permutation-based transformation for key matrix for illustration purposes.

Auxiliary’s backpropagation through linear layer For a linear layer as defined in Definition D.1, the linear backpropagation layer takes in the loss gradient w.r.t. output (∂y) and computes the loss gradient w.r.t. input (∂x).

Definition D.2 (Linear backpropagation). For a weight $\mathbf{W} \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$, the linear backpropagation layer takes $\partial y \in \mathbb{R}^{D_{\text{aux}}}$ as input and outputs $\partial x = \mathbf{W}^\top \partial y$.

TINT Linear backpropagation module This module will aim to simulate the auxiliary’s linear backpropagation. The input embedding e_t to this module will contain the gradient of the loss w.r.t. y_t , i.e. ∂y_t . As given in Definition D.2, this module will output the gradient of the loss w.r.t. x_t , given by $\partial x_t = \mathbf{W}^\top \partial y_t$.

We first use the residual connection to copy the prefix embeddings $\{v_j\}$ (i.e., the rows of \mathbf{W}) from the forward propagation module. A straightforward construction would be to use the Linear Forward module but with the columns of \mathbf{W} stored in the prefix tokens, thereby simulating multiplication with \mathbf{W}^\top . However, such a construction requires applying attention to the prefix tokens, which increases the size of the construction substantially.

We instead perform the operation more efficiently by splitting it across attention heads. In particular, once we view the operation as $\partial x_t = \sum_i (\partial y_t)_i w_i$, we can see that the attention score between the current token and the prefix token containing w_i must be $(\partial y_t)_i$. Using value vectors as rows of \mathbf{W} returns the desired output. Similar to the Linear Forward module, we shard the weights into S' parts to parallelize across more attention heads. Please see Figure 3.

Auxiliary’s linear descent update Finally, the linear descent layer updates the weight and the bias parameters using a batch of inputs $\{x_t\}_{t \leq T_{\text{aux}}}$ and the loss gradient w.r.t. the corresponding outputs $\{\partial y_t\}_{t \leq T_{\text{aux}}}$.

Definition D.3 (Linear descent). For a weight $\mathbf{W} \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$ and a bias $\mathbf{b} \in \mathbb{R}^{D_{\text{aux}}}$, the linear descent layer takes in a batch of inputs $\{x_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and gradients $\{\partial y_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and updates the parameters as follows:

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \sum_{t \leq T_{\text{aux}}} \partial y_t x_t^\top; \quad \mathbf{b} \leftarrow \mathbf{b} - \eta \sum_{t \leq T_{\text{aux}}} \partial y_t.$$

TINT Linear descent module The input embedding e_t to this module will contain the gradient of the loss w.r.t. y_t , i.e. ∂y_t .

As in the Linear backpropagation module, the prefix tokens $\{v_j\}$ will contain the rows of \mathbf{W} and \mathbf{b} , which have been copied from the Linear forward module using residual connections. Since, in

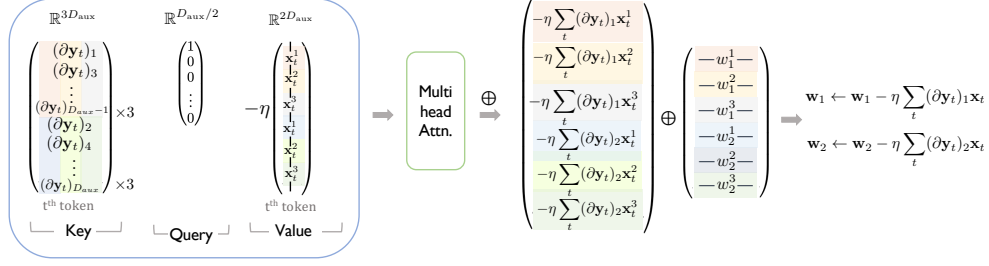


Figure 4: TINT computes the parameter gradients for a linear layer as a H -head attention layer ($H = 6$ pictured), with the gradient of the loss w.r.t. linear layer output (∂y_t) as the query, the positional one-hot vector of prefix embeddings as the key, and the input to the linear layer (x_t) as the value. The auxiliary model parameters in the prefix embeddings are then updated using a residual connection. Similar to the Linear Forward module (Figure 2), we distribute the dot product computations across all attention heads, by sharding the vectors into S' ($S' = 3$ here) parts. We omitted the identical transformation for query, and value matrices, and permutation-based transformation for key matrix for simplicity.

addition to the gradients, we also require the input to the linear layer, we will use residual connections to copy the input $\{x_t\}$ to their respective embeddings $\{e_t\}$, from the Linear Forward module. As given in Definition D.3, this module will update \mathbf{W} and \mathbf{b} using the gradient descent rule.

Focusing on w_i , the descent update is given by $w_i \leftarrow w_i - \eta \sum_t (\partial y_t)_i x_t$. For the prefix token v_j that contains w_i , the update term $-\eta \sum_t (\partial y_t)_i x_t$ can be expressed with an attention head that represents the attention between the prefix token v_j and any token e_t with score $(\partial y_t)_i$ and value $-\eta x_t$. The residual connection can then be used to update the weights w_i in v_j .

For the bias \mathbf{b} , the descent update is given by $\mathbf{b} \leftarrow \mathbf{b} - \eta \sum_t \partial y_t$. With \mathbf{b} present in v_1 , we use one attention head to represent the attention score between prefix token v_1 and any token e_t as 1, with the value being $-\eta \partial y_t$. The residual connection can then be used to update the weights \mathbf{b} in v_1 .

The above process can be further parallelized across multiple attention heads, by sharding each weight computation into S' parts. Please see Figure 4.

D.1 H_{sim} -split operation

We leverage local structure within the linear operations of TINT to make the construction smaller. We build two H_{sim} -split operations to replace all the linear operations. We use d_{sim} to denote $D_{\text{sim}}/H_{\text{sim}}$ in the following definitions.

Definition D.4 (Split-wise H_{sim} -split Linear operation). For weight and bias parameters $\mathbf{W}^{\text{TINT}} \in \mathbb{R}^{H_{\text{sim}} \times d_{\text{sim}} \times d_{\text{sim}}}$, $\mathbf{B}^{\text{TINT}} \in \mathbb{R}^{H_{\text{sim}} \times d_{\text{sim}}}$, this layer takes in input $e \in \mathbb{R}^{D_{\text{sim}}}$ and returns $\tilde{e} = \text{VECTORIZE}(\tilde{\mathbf{S}} + \mathbf{B}^{\text{TINT}})$, with $\tilde{\mathbf{S}} \in \mathbb{R}^{H_{\text{sim}} \times d_{\text{sim}}}$ defined with rows $\{\mathbf{W}_h^{\text{TINT}} \text{SPLIT}_{H_{\text{sim}}}(e)_h\}_{h \leq H_{\text{sim}}}$.

Definition D.5 (Dimension-wise H_{sim} -split Linear operation). For weight and bias parameters $\mathbf{W}^{\text{TINT}} \in \mathbb{R}^{d_{\text{sim}} \times H_{\text{sim}} \times H_{\text{sim}}}$, $\mathbf{B}^{\text{TINT}} \in \mathbb{R}^{d_{\text{sim}} \times H_{\text{sim}}}$, this layer takes in input $e \in \mathbb{R}^{D_{\text{sim}}}$, defines $\mathbf{S} \in \mathbb{R}^{d_{\text{sim}} \times H_{\text{sim}}}$ with columns $\{\text{SPLIT}_{H_{\text{sim}}}(e)_h\}_{h \leq H_{\text{sim}}}$, and returns $\tilde{e} = \text{VECTORIZE}((\tilde{\mathbf{S}} + \mathbf{B}^{\text{TINT}})^\top)$, where $\tilde{\mathbf{S}} \in \mathbb{R}^{d_{\text{sim}} \times H_{\text{sim}}}$ is defined with rows $\{\mathbf{W}_d^{\text{TINT}} \mathbf{s}_d^{\text{TINT}}\}_{d \leq d_{\text{sim}}}$.

We find that we can replace all the linear operations with a splitwise H_{sim} -split Linear operation followed by a dimensionwise H_{sim} -split Linear operation, and an additional splitwise H_{sim} -split Linear operation, if necessary. A linear operation on D_{sim} -dimensional space involves D_{sim}^2 parameters, while its replacement requires $D_{\text{sim}}^2/H_{\text{sim}} + 2D_{\text{sim}}H_{\text{sim}}$ parameters, effectively reducing the total number of necessary parameters by H_{sim} .

We motivate the H_{sim} -split linear operations with an example. We consider the Linear Forward module in Figure 2 for simulating a linear operation with parameters $\mathbf{W} \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$ and no biases. For simplicity of presentation, we assume D_{aux} is divisible by 4. We stack 2 rows of weights per

prefix embedding. We distribute the dot-product computation across the $H_{\text{sim}} = 6$ attention heads, by sharding each weight into 3 parts. Since we require to have enough space to store all the sharded computation from the linear attention heads, we require $D_{\text{sim}} = 3D_{\text{aux}}$ (we get 3 values for each of the D_{aux} weights in \mathbf{W}). For presentation, for a given vector $\mathbf{v} \in \mathbb{R}^{D_{\text{aux}}}$, we represent $\text{SPLIT}_3(\mathbf{v})_i$ by \mathbf{v}^i for all $1 \leq i \leq 3$.

Now, consider the final linear operation responsible for combining the output of the attention heads. The output, after the linear operation, should contain $\mathbf{W}\mathbf{x}_t$ in the first D_{aux} coordinates. At any position t , if we stack the output of the linear attention heads as rows of a matrix $\mathbf{S}_t \in \mathbb{R}^{H_{\text{sim}} \times D_{\text{sim}}/H_{\text{sim}}}$ we get

$$\mathbf{S}_t = \begin{bmatrix} \langle \mathbf{w}_1^1, \mathbf{x}_t^1 \rangle & \langle \mathbf{w}_3^1, \mathbf{x}_t^1 \rangle & \langle \mathbf{w}_5^1, \mathbf{x}_t^1 \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}-1}^1, \mathbf{x}_t^1 \rangle \\ \langle \mathbf{w}_1^2, \mathbf{x}_t^2 \rangle & \langle \mathbf{w}_3^2, \mathbf{x}_t^2 \rangle & \langle \mathbf{w}_5^2, \mathbf{x}_t^2 \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}-1}^2, \mathbf{x}_t^2 \rangle \\ \langle \mathbf{w}_1^3, \mathbf{x}_t^3 \rangle & \langle \mathbf{w}_3^3, \mathbf{x}_t^3 \rangle & \langle \mathbf{w}_5^3, \mathbf{x}_t^3 \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}-1}^3, \mathbf{x}_t^3 \rangle \\ \langle \mathbf{w}_2^1, \mathbf{x}_t^1 \rangle & \langle \mathbf{w}_4^1, \mathbf{x}_t^1 \rangle & \langle \mathbf{w}_6^1, \mathbf{x}_t^1 \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}}^1, \mathbf{x}_t^1 \rangle \\ \langle \mathbf{w}_2^2, \mathbf{x}_t^2 \rangle & \langle \mathbf{w}_4^2, \mathbf{x}_t^2 \rangle & \langle \mathbf{w}_6^2, \mathbf{x}_t^2 \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}}^2, \mathbf{x}_t^2 \rangle \\ \langle \mathbf{w}_2^3, \mathbf{x}_t^3 \rangle & \langle \mathbf{w}_4^3, \mathbf{x}_t^3 \rangle & \langle \mathbf{w}_6^3, \mathbf{x}_t^3 \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}}^3, \mathbf{x}_t^3 \rangle \end{bmatrix}$$

Note that for each $j \leq D_{\text{aux}}$, we have $\langle \mathbf{w}_j, \mathbf{x}_t \rangle = \sum_{i=1}^3 \langle \mathbf{w}_j^i, \mathbf{x}_t^i \rangle$. Thus, with a column-wise linear operation on \mathbf{S}_t , we can sum the relevant elements in each column to get

$$\mathbf{S}_t^{\text{col}} = \begin{bmatrix} \langle \mathbf{w}_1, \mathbf{x}_t \rangle & \langle \mathbf{w}_3, \mathbf{x}_t \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}/2-1}, \mathbf{x}_t \rangle & 0 & 0 & \cdots & 0 \\ \langle \mathbf{w}_2, \mathbf{x}_t \rangle & \langle \mathbf{w}_4, \mathbf{x}_t \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}/2}, \mathbf{x}_t \rangle & 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & \langle \mathbf{w}_{D_{\text{aux}}/2+1}, \mathbf{x}_t \rangle & \langle \mathbf{w}_{D_{\text{aux}}/2+3}, \mathbf{x}_t \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}-1}, \mathbf{x}_t \rangle \\ 0 & 0 & \cdots & 0 & \langle \mathbf{w}_{D_{\text{aux}}/2+2}, \mathbf{x}_t \rangle & \langle \mathbf{w}_{D_{\text{aux}}/2+4}, \mathbf{x}_t \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}}, \mathbf{x}_t \rangle \\ 0 & 0 & \cdots & 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 0 & 0 & \cdots & 0 \end{bmatrix}$$

A row-wise linear operation on $\mathbf{S}_t^{\text{col}}$ can space out the non-zero elements in the matrix and give us

$$\mathbf{S}_t^{\text{row}} = \begin{bmatrix} \langle \mathbf{w}_1, \mathbf{x}_t \rangle & 0 & \langle \mathbf{w}_3, \mathbf{x}_t \rangle & 0 & \cdots & \langle \mathbf{w}_{D_{\text{aux}}/2-1}, \mathbf{x}_t \rangle & 0 \\ 0 & \langle \mathbf{w}_2, \mathbf{x}_t \rangle & 0 & \langle \mathbf{w}_4, \mathbf{x}_t \rangle & \cdots & 0 & \langle \mathbf{w}_{D_{\text{aux}}/2}, \mathbf{x}_t \rangle \\ \langle \mathbf{w}_{D_{\text{aux}}/2+1}, \mathbf{x}_t \rangle & 0 & \langle \mathbf{w}_{D_{\text{aux}}/2+3}, \mathbf{x}_t \rangle & 0 & \cdots & \langle \mathbf{w}_{D_{\text{aux}}-1}, \mathbf{x}_t \rangle & 0 \\ 0 & \langle \mathbf{w}_{D_{\text{aux}}/2+2}, \mathbf{x}_t \rangle & 0 & \langle \mathbf{w}_{D_{\text{aux}}/2+4}, \mathbf{x}_t \rangle & \cdots & 0 & \langle \mathbf{w}_{D_{\text{aux}}}, \mathbf{x}_t \rangle \\ 0 & 0 & \cdots & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & \cdots & 0 & 0 \end{bmatrix}$$

Finally, a column-wise linear operation on $\mathbf{S}_t^{\text{row}}$ helps to get the non-zero elements in the correct order.

$$\bar{\mathbf{S}}_t^{\text{col}} = \begin{bmatrix} \langle \mathbf{w}_1, \mathbf{x}_t \rangle & \langle \mathbf{w}_2, \mathbf{x}_t \rangle & \langle \mathbf{w}_3, \mathbf{x}_t \rangle & \langle \mathbf{w}_4, \mathbf{x}_t \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}/2-1}, \mathbf{x}_t \rangle & \langle \mathbf{w}_{D_{\text{aux}}/2}, \mathbf{x}_t \rangle \\ \langle \mathbf{w}_{D_{\text{aux}}/2+1}, \mathbf{x}_t \rangle & \langle \mathbf{w}_{D_{\text{aux}}/2+2}, \mathbf{x}_t \rangle & \langle \mathbf{w}_{D_{\text{aux}}/2+3}, \mathbf{x}_t \rangle & \langle \mathbf{w}_{D_{\text{aux}}/2+4}, \mathbf{x}_t \rangle & \cdots & \langle \mathbf{w}_{D_{\text{aux}}-1}, \mathbf{x}_t \rangle & \langle \mathbf{w}_{D_{\text{aux}}}, \mathbf{x}_t \rangle \\ 0 & 0 & 0 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \cdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 0 & 0 \end{bmatrix}$$

The desired output is then given by $\text{VECTORIZE}(\{\bar{\mathbf{s}}_{t,j}^{\text{col}}\}_{j=1}^{D_{\text{aux}}})$, which contains $\mathbf{W}\mathbf{x}_t$ in the first D_{aux} coordinates. The operations that convert \mathbf{S}_t to $\mathbf{S}_t^{\text{col}}$ and $\mathbf{S}_t^{\text{row}}$ to $\bar{\mathbf{S}}_t^{\text{row}}$ represents a split-wise 6-split linear operation, while the operation that converts $\mathbf{S}_t^{\text{col}}$ to $\mathbf{S}_t^{\text{row}}$ represents a dimension-wise 6-split linear operation. A naive linear operation on the output of the attention heads would require D_{sim}^2 parameters, while its replacement requires $D_{\text{sim}}^2/6$ parameters to represent a dimension-wise 6-split linear operation, and an additional $12D_{\text{sim}}$ parameters to represent the split-wise 6-split linear operations.

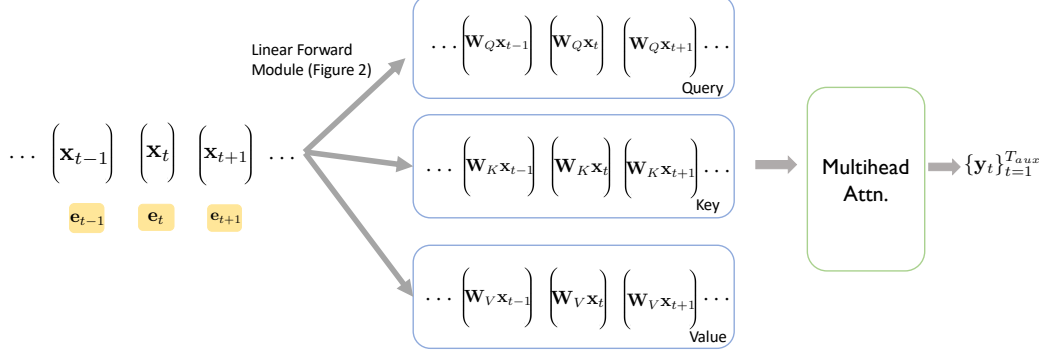


Figure 5: TINT simulates the forward pass of a self-attention layer of the auxiliary model with a Linear Forward module (Figure 2) and a TINT softmax attention layer (Definition C.1). The Linear Forward module computes the query, key, and value vectors using a Linear Forward module on the current embeddings, changing the prefix embeddings to correspond to W_Q , W_K , and W_V respectively.

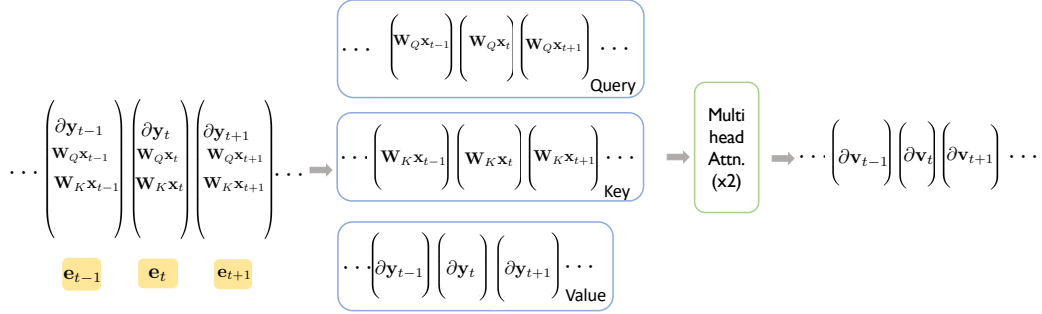


Figure 6: The gradient w.r.t. the value vectors $\{\partial v_t\}$ (Definition E.2) forms the integral component for both TINT self-attention backward and descent update modules. TINT computes $\{\partial v_t\}$ using a softmax attention and a linear attention layer. We first use residual connections to copy the query and key vectors to the current embeddings from the TINT Self-attention Forward module (Figure 5). The softmax attention layer re-computes the attention scores $\{a_{t,j}^h\}$ between all token pairs $\{(t, j)\}$ and stores them in the token embeddings. The linear attention layer uses the one-hot position embeddings of the input tokens as the query to use the transposed attention scores $\{a_{j,t}^h\}$ for all token pairs $\{(t, j)\}$ and use the gradients $\{\partial y_t\}$ as the value vectors to compute $\{\partial v_t\}$.

E Self-attention layer

We first introduce multi-head attention, generalizing single-head attention (Definition A.5).

Definition E.1 (Auxiliary self-attention with H_{aux} heads). For query, key, and value weights $W_Q, W_K, W_V \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$ and bias $b_Q, b_K, b_V \in \mathbb{R}^{D_{\text{aux}}}$, a self-attention layer with H_{aux} attention heads and a function $f_{\text{attn}} : \mathbb{R}^{T_{\text{aux}}} \rightarrow \mathbb{R}^{T_{\text{aux}}}$ takes a sequence $\{x_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ as input and outputs $\{y_t\}_{t \leq T_{\text{aux}}}$, with

$$y_t = \text{VECTORIZE}\left(\sum_{j \leq T_{\text{aux}}} a_{t,j}^h v_j^h\right)_{h \leq H_{\text{aux}}}. \quad (6)$$

$a_{t,j}^h$ is defined as the attention score of head h between tokens at positions t and j , and is given by

$$a_{t,j}^h = \text{softmax}(\mathbf{K}^h \mathbf{q}_t^h)_j. \quad (7)$$

Here, q_t, k_t, v_t denote the query, key, and value vectors at each position t , computed as $W_Q x_t + b_Q$, $W_K x_t + b_K$, and $W_V x_t + b_V$ respectively. In addition, q_t^h, k_t^h, v_t^h denote $\text{SPLIT}_{H_{\text{aux}}}(\mathbf{q}_t)_h$, $\text{SPLIT}_{H_{\text{aux}}}(\mathbf{k}_t)_h$, and $\text{SPLIT}_{H_{\text{aux}}}(\mathbf{v}_t)_h$ respectively for all $t \leq T_{\text{aux}}$, and $h \leq H_{\text{aux}}$. $\mathbf{K}^h \in \mathbb{R}^{T_{\text{aux}} \times D_{\text{aux}}}$ is defined with its rows as $\{k_t^h\}_{t \leq T_{\text{aux}}}$ for all $h \leq H_{\text{aux}}$.

In the discussions below, we consider a self-attention layer in the auxiliary model with parameters $\{\mathbf{W}_Q, \mathbf{b}_Q, \mathbf{W}_K, \mathbf{b}_K, \mathbf{W}_V, \mathbf{b}_V\}$ that takes in input sequence $\mathbf{x}_1, \dots, \mathbf{x}_{T_{\text{aux}}}$ and outputs $\mathbf{y}_1, \dots, \mathbf{y}_{T_{\text{aux}}}$, with $\{\mathbf{y}_t\}_{t=1}^{T_{\text{aux}}}$ given by (6). As in the definition, $\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t$ denote the query, key, and value vectors for position t . We will use TINT self-attention modules in order to simulate the operations on the auxiliary’s self-attention layer. To do so, we will need $H_{\text{sim}} \geq H_{\text{aux}}$ in the corresponding TINT self-attention modules.

TINT Self-attention forward module The input embedding to this module \mathbf{e}_t at each position t will contain \mathbf{x}_t in its first D_{aux} coordinates. The self-attention module can be divided into four sub-operations: Computation of (a) query vectors $\{\mathbf{q}_t\}_{t \leq T}$, (b) key vectors $\{\mathbf{k}_t\}_{t \leq T}$, (c) value vectors $\{\mathbf{v}_t\}_{t \leq T}$, and (d) $\{\mathbf{y}_t\}_{t \leq T}$ using (6). Please see Figure 5.

- Sub-operations (a): The computation of query vector $\mathbf{q}_t := \mathbf{W}_Q \mathbf{x}_t + \mathbf{b}_Q$ at each position t is a linear operation involving parameters $\mathbf{W}_Q, \mathbf{b}_Q$. Thus, we can first feed in the stacked rows of \mathbf{W}_Q and \mathbf{b}_Q onto the prefix embeddings $\{\mathbf{v}_j\}$. We use a Linear Forward module (Appendix D) on the current embeddings and the prefix embeddings to get embedding \mathbf{e}_t^q at each position t that contains \mathbf{q}_t in the first D_{aux} coordinates.
- Sub-operations (b, c): Similar to (a), we feed in the stacked rows of the necessary parameters onto the prefix embeddings $\{\mathbf{v}_j\}$, and call two Linear Forward Modules (Appendix D) independently to get embeddings \mathbf{e}_t^k , and \mathbf{e}_t^v containing \mathbf{k}_t and \mathbf{v}_t respectively.
We now combine the embeddings $\mathbf{e}_t^q, \mathbf{e}_t^k$, and \mathbf{e}_t^v to get an embedding \mathbf{e}_t that contain $\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t$ in the first $3D_{\text{aux}}$ coordinates.
- Sub-operation (d): Finally, we call a TINT self-attention module (Definition C.1) on our current embeddings $\{\mathbf{e}_t\}_{t \leq T}$ to compute $\{\mathbf{y}_t\}_{t \leq T}$. The query, key, and value parameters in the self-attention module contain sub-Identity blocks that pick out the relevant information from $\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t$ stored in \mathbf{e}_t .

Remark: Sub-operations (a), (b), and (c) can be represented as a single linear operation with a weight $\mathbf{W} \in \mathbb{R}^{3D_{\text{aux}} \times D_{\text{aux}}}$ by concatenating the rows of $\{\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V\}$ and a bias $\mathbf{b} \in \mathbb{R}^{3D_{\text{aux}}}$ that concatenates $\{\mathbf{b}_Q, \mathbf{b}_K, \mathbf{b}_V\}$. Thus, they can be simulated with a single Linear Forward Module, with \mathbf{W}, \mathbf{b} fed into the prefix embeddings. However, we decide to separate them in order to limit the number of prefix embeddings and the embedding size. E.g. for GPT-2, $D_{\text{aux}} = 768$. This demands either a $3 \times$ increase in the embedding size in TINT or a $3 \times$ increase in the number of prefix embeddings. Hence, in order to minimize the parameter cost, we call Linear Forward Module separately to compute $\mathbf{q}_t, \mathbf{k}_t$, and \mathbf{v}_t at each position t .

Auxiliary’s backpropagation through self-attention For an auxiliary self-attention layer as defined in Definition E.1, the backpropagation layer takes in the loss gradient w.r.t. output ($\{\partial_{\mathbf{y}_t}\}_{t \leq T_{\text{aux}}}$) and computes the loss gradient w.r.t. input token ($\{\partial_{\mathbf{x}_t}\}_{t \leq T_{\text{aux}}}$).

Definition E.2. [Auxiliary self-attention backpropagation] For query, key, and value weights $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$ and bias $\mathbf{b}_Q, \mathbf{b}_K, \mathbf{b}_V \in \mathbb{R}^{D_{\text{aux}}}$, the backpropagation layer corresponding to a self-attention layer with H_{aux} attention heads takes a sequence $\{\partial_{\mathbf{y}_t} \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and $\{\mathbf{x}_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ as input and outputs $\{\partial_{\mathbf{x}_t}\}_{t \leq T_{\text{aux}}}$, with

$$\begin{aligned} \partial_{\mathbf{x}_t} &= \mathbf{W}_Q^\top \partial_{\mathbf{q}_t} + \mathbf{W}_K^\top \partial_{\mathbf{k}_t} + \mathbf{W}_V^\top \partial_{\mathbf{v}_t}, \quad \text{with} \\ \partial_{\mathbf{q}_t} &= \text{VECTORIZE}(\{\sum_j a_{t,j}^h ((\partial_{\mathbf{y}_t}^h)^\top \mathbf{v}_j^h) [\mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h]\}_{h \leq H_{\text{aux}}}); \\ \partial_{\mathbf{k}_t} &= \text{VECTORIZE}(\{\sum_j a_{j,t}^h \mathbf{q}_j^h [(\partial_{\mathbf{y}_j}^h)^\top (\mathbf{v}_t^h - \sum_{j'} a_{j,j'}^h \mathbf{v}_{j'}^h)]\}_{h \leq H_{\text{aux}}}); \\ \partial_{\mathbf{v}_t} &= \text{VECTORIZE}(\{\sum_j a_{j,t}^h \partial_{\mathbf{y}_j}^h\}_{h \leq H_{\text{aux}}}) \end{aligned}$$

Here, $\mathbf{q}_t, \mathbf{k}_t$, and \mathbf{v}_t refer to query, key, and value vectors at each position t , with the attention scores $\{a_{t,j}^h\}_{t,j \leq T_{\text{aux}}, h \leq H_{\text{aux}}}$.

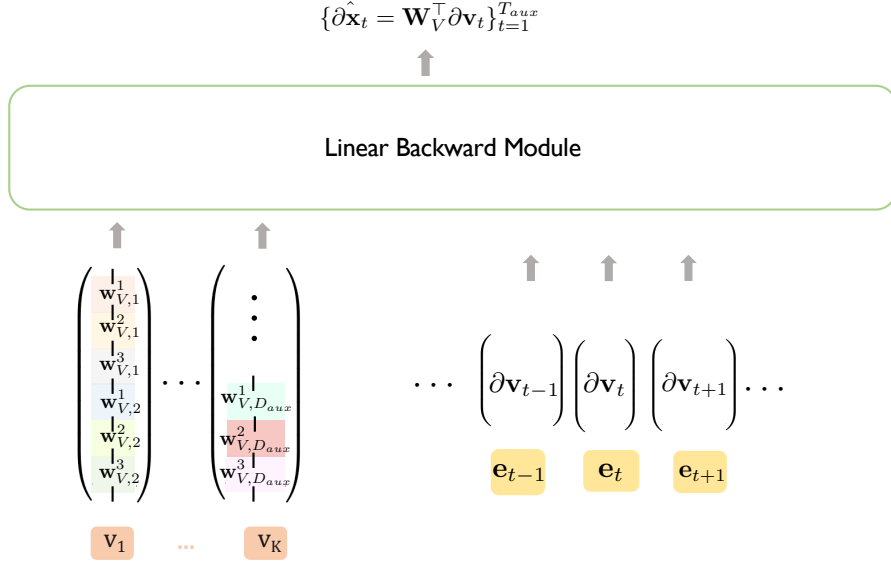


Figure 7: TINT simulates the backward pass of a self-attention layer of the auxiliary model using a Linear Backward module (Figure 3). The input embeddings contain the gradient of the loss w.r.t. the value vectors (∂_{v_t}) computed in Figure 6. The value matrix \mathbf{W}_V is encoded in the prefix embeddings. We call the Linear Backward module on this sequence.

Complexity of true backpropagation The much-involved computation in the above operation is due to the computation of ∂_{q_t} and ∂_{k_t} at each position t . For the following discussion, we assume that our current embeddings e_t contain q_t, k_t, v_t , in addition to the gradient ∂_{y_t} . The computation of ∂_{q_t} (and similarly ∂_{k_t}) at any position t involves the following sequential computations and the necessary TINT modules.

- $\{(\partial_{y_t^h})^\top v_j^h\}_{j \leq T_{\text{aux}}}\}_{h \leq H_{\text{aux}}}$ with a TINT linear self-attention module (Definition C.1), with atleast H_{aux} attention heads that represent the attention score between e_t and any other token e_j , by $\{(\partial_{y_t^h})^\top v_j^h\}_{h \leq H_{\text{aux}}}$.
- Attention scores $\{a_{t,j}^h\}_{h \leq H_{\text{aux}}}$, which requires a TINT softmax self-attention module (Definition C.1), with atleast H_{aux} heads, that uses the already present $\{q_t, k_t, v_t\}$ in the current embeddings e_t to re-compute the attention scores.
- $\{a_{t,j}^h (\partial_{y_t^h})^\top v_j^h\}_{h \leq H_{\text{aux}}}$ for all $j \leq T_{\text{aux}}$ by multiplying the attention scores $\{a_{t,j}^h\}_{h \leq H_{\text{aux}}}$ with $\{(\partial_{y_t^h})^\top v_j^h\}_{h \leq H_{\text{aux}}}$ using an MLP layer (Lemma C.4). Furthermore, $\{\sum_j a_{t,j}^h k_j^h\}_{h \leq H_{\text{aux}}}$ needs to be computed in parallel as well, with additional attention heads.
- ∂_{y_t} with a TINT linear self-attention module (Definition C.1), with atleast H_{aux} attention heads that represent the attention score between any token e_j and e_t by $\{a_{t,j}^h (\partial_{y_t^h})^\top v_j^h\}_{h \leq H_{\text{aux}}}$, with value vectors given by $\{k_j^h - \sum_{j'} a_{t,j'}^h k_{j'}^h\}_{h \leq H_{\text{aux}}}$.

The sequential computation requires the simulator to store $\{(\partial_{y_t^h})^\top v_j^h\}_{j \leq T_{\text{aux}}}\}_{h \leq H_{\text{aux}}}$ and $\{a_{t,j}^h\}_{h \leq H_{\text{aux}}}$ in the token embedding e_t , which requires an additional $2T_{\text{aux}}H_{\text{aux}}$ embedding dimension size. To avoid the much-involved computation for the true gradient propagation, we instead only use the gradients w.r.t. v_t .

Approximate auxiliary self-attention backpropagation We formally extend the definition of approximate gradients $\{\partial_{x_t}\}_{t=1}^{T_{\text{aux}}}$ from Definition A.7 to multi-head attention in Definition E.6.

In the upcoming theorem, we formally show that if on a given sequence $\{x_t\}_{t \leq T_{\text{aux}}}$, for all token positions all the attention heads in a self-attention layer primarily attend to a single token, then the approximate gradient $\widehat{\partial_{x_t}}$ is close to the true gradient ∂_{x_t} at each position t .

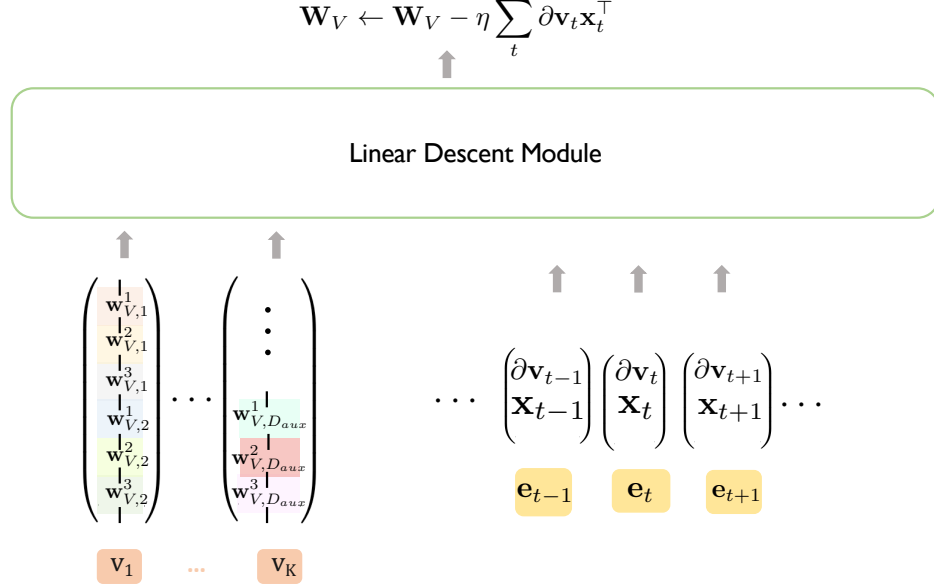


Figure 8: TINT simulates the backward pass of the self-attention layer in the auxiliary model by employing the Linear Descent module (Figure 4). The input embeddings consist of the gradient of the loss with respect to the value vectors (∂_{v_t}) computed in Figure 6. Additionally, we incorporate a residual connection to copy the input from the Self-attention Forward module (Figure 5) into \mathbf{x}_t . Before invoking the Linear Descent module, we represent the value parameters (\mathbf{W}_V) into the prefix embeddings. TINT simulates the backward pass of a self-attention layer of the auxiliary model using a Linear Descent module (Figure 4).

Definition E.3 (ε -hard attention head). For the Self-Attention layer of H_{aux} heads in Definition E.1, on a given input sequence $\{\mathbf{x}_t\}_{t=1}^{T_{\text{aux}}}$, an attention head $h \leq H_{\text{aux}}$ is defined to be ε -hard on the input sequence, if for all positions $t \leq T_{\text{aux}}$, there exists a position $t_0 \leq T_{\text{aux}}$ such that $a_{t,t_0}^h \geq 1 - \varepsilon$.

Theorem E.4. With the notations in Definitions A.7, E.1 and E.2, if on a given input sequence $\{\mathbf{x}_t\}_{t=1}^{T_{\text{aux}}}$, with its query, key, and value vectors $\{\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t\}_{t=1}^{T_{\text{aux}}}$, all the H_{aux} attention heads are ε -hard for some $\varepsilon > 0$, then for a given sequence of gradients $\{\partial_{y_t}\}_{t=1}^{T_{\text{aux}}}$,

$$\|\partial_{\mathbf{q}_t}\|_2, \|\partial_{\mathbf{k}_t}\|_2 \leq \mathcal{O}(\varepsilon B_x^2 B_w^2 B_y), \quad \text{for all } t \leq T_{\text{aux}},$$

where $B_x = \max_{t \leq T_{\text{aux}}} \|\mathbf{x}_t\|_2$, $B_y = \max_{t \leq T_{\text{aux}}} \|\partial_{y_t}\|_2$, and $B_w = \max\{\|\mathbf{W}_K\|_2, \|\mathbf{W}_Q\|_2, \|\mathbf{W}_V\|_2, \|\mathbf{b}_K\|_2, \|\mathbf{b}_Q\|_2, \|\mathbf{b}_V\|_2\}$.

This implies, for each position t , $\|\widehat{\partial_{\mathbf{x}_t}} - \partial_{\mathbf{x}_t}\|_2 \leq \mathcal{O}(\varepsilon B_x^2 B_w^3 B_y)$.

TINT Self-attention backpropagation module The input embeddings \mathbf{e}_t contain ∂_{y_t} in the first D_{aux} coordinates. Since we require to re-compute the attention scores $\{a_{t,j}^h\}_{j \leq T_{\text{aux}}, h \leq H_{\text{aux}}}$, we need to copy the query, key, and value vectors $\mathbf{q}_t, \mathbf{k}_t$, and \mathbf{v}_t from the TINT self-attention Forward module at each position t . Furthermore, we use the residual connection to copy the prefix embeddings $\{\mathbf{v}_j\}$, which contain the rows of \mathbf{W}_V , from the TINT self-attention Forward module.

The operation can be divided into three sub-operations: Computing (a) attention scores $\{a_{t,j}^h\}_{h \leq H_{\text{aux}}}$ for all $j \leq T_{\text{aux}}$, at each position t , (b) ∂_{v_t} from $\{a_{t,j}^h\}_{h \leq H_{\text{aux}}}$ and ∂_{y_t} , and (c) $\widehat{\partial_{\mathbf{x}_t}}$ from ∂_{v_t} .

- Sub-operation (a): Since, the current embeddings \mathbf{e}_t contain $\mathbf{q}_t, \mathbf{k}_t$, we can simply call a self-attention attention module to compute the attention scores $\{a_{t,j}^h\}_{h \leq H_{\text{aux}}}$ for all $j \leq T$ and store them in the current embeddings. We further retain ∂_{y_t} and \mathbf{v}_t for further operations using residual connections.

- Sub-operation (b): With the current embeddings e_t containing the attention scores $\{a_{t,j}^h\}_{h \leq H_{\text{aux}}}$ for all $j \leq T$, and the gradient ∂_{y_t} , we can compute ∂_{v_t} using a TINT linear self-attention module with atleast H_{aux} attention heads, that represent the attention scores between tokens e_t and e_j for any j as $\{a_{j,t}^h\}_{h \leq H_{\text{aux}}}$ and use $\text{SPLIT}_{H_{\text{aux}}}(\partial_{y_t})$ as their value vectors.
- Sub-operation (c): And finally, the computation of $\widehat{\partial_{x_t}}$ is identical to the backpropagation through a linear layer, with parameters \mathbf{W}_V and \mathbf{b}_V . Hence, we call a Linear backpropagation module on the current embeddings, that contain ∂_{y_t} and the prefix embeddings that contain \mathbf{W}_V and \mathbf{b}_V .

Separating sub-operations (a) and (b) The operation for computing ∂_{v_t} in Definition A.7 looks very similar to the computation of y_t in Equation (6). However, the major difference is that instead of the attention scores being $\{a_{t,j}^h\}_{h \leq H_{\text{aux}}}$ between token t and any token j , we need the attention scores to be $\{a_{j,t}^h\}_{h \leq H_{\text{aux}}}$. Thus, unless our model allows a transpose operation on the attention scores, we need to first store them in our embeddings and then use an additional self-attention module that can pick the right attention scores between tokens using position embeddings. Please see Figure 7.

Auxiliary’s value descent update Similar to the complexity of true backpropagation, the descent updates for $\mathbf{W}_Q, \mathbf{b}_Q, \mathbf{W}_K, \mathbf{b}_K$ are quite expensive to express with the transformer layers. Hence, we focus simply on updating on $\mathbf{W}_V, \mathbf{b}_V$, while keeping the others fixed.

Definition E.5 (Auxiliary self-attention value descent). For query, key, and value weights $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$ and bias $\mathbf{b}_Q, \mathbf{b}_K, \mathbf{b}_V \in \mathbb{R}^{D_{\text{aux}}}$, the value descent layer corresponding to a self-attention layer with H_{aux} attention heads and any function $f_{\text{attn}} : \mathbb{R}^{T_{\text{aux}}} \rightarrow \mathbb{R}^{T_{\text{aux}}}$ takes in a batch of gradients $\{\partial_{y_t} \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and inputs $\{x_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and updates $\mathbf{W}_V, \mathbf{b}_V$ as follows:

$$\mathbf{W}_V \leftarrow \mathbf{W}_V - \eta \sum_{t \leq T_{\text{aux}}} \partial_{v_t} x_t^\top, \quad \mathbf{b}_V \leftarrow \mathbf{b}_V - \eta \sum_{t \leq T_{\text{aux}}} \partial_{v_t},$$

$$\text{where } \partial_{v_t} = \text{VECTORIZE}\left(\left\{\sum_j a_{j,t}^h \partial_{y_j^h}\right\}_{h \leq H_{\text{aux}}}\right)$$

Here, v_t refers to value vectors at each position t , as defined in Definition E.1.

TINT Self-attention descent module The input embeddings contain ∂_{v_t} in the first D_{aux} coordinates, from the TINT self-attention backpropagation module. Furthermore, the prefix embeddings $\{v_j\}$ contain the stacked rows of \mathbf{W}_V and \mathbf{b}_V , continuing from the TINT self-attention backpropagation module.

Since we further need the input x_t to the auxiliary self-attention layer under consideration, we use residual connections to copy x_t from the TINT self-attention Forward module at each position t .

The updates of \mathbf{W}_V and \mathbf{b}_V are equivalent to the parameter update in a linear layer, involving gradients $\{\partial_{v_t}\}$ and input $\{x_t\}$. Thus, we call a Linear descent module on the current embeddings and the prefix embeddings to get the updated value parameters. Please see Figure 8.

E.1 Approximate auxiliary self-attention backpropagation

Definition E.6. For query, key, and value weights $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$ and bias $\mathbf{b}_Q, \mathbf{b}_K, \mathbf{b}_V \in \mathbb{R}^{D_{\text{aux}}}$, the approximate backpropagation layer corresponding to a self-attention layer with H_{aux} attention heads takes a sequence $\{\partial_{y_t} \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and $\{x_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ as input and outputs $\{\partial_{x_t} := \text{VECTORIZE}(\{\partial_{x_t^h}\}_{h \leq H_{\text{aux}}})\}_{t \leq T_{\text{aux}}}$, with

$$\widehat{\partial_{x_t}} = \mathbf{W}_V^\top \partial_{v_t}, \quad \text{where } \partial_{v_t} = \text{VECTORIZE}\left(\left\{\sum_j a_{j,t}^h \partial_{y_j^h}\right\}_{h \leq H_{\text{aux}}}\right)$$

Here, q_t, k_t , and v_t refer to query, key, and value vectors at each position t , as defined in Definition E.1, with the attention scores $\{a_{t,j}^h\}_{t,j \leq T_{\text{aux}}, h \leq H_{\text{aux}}}$ defined in Equation (7).

E.2 Proofs of theorems and gradient definitions

We restate the theorems and definitions, before presenting their proofs for easy referencing.

Definition E.2. [Auxiliary self-attention backpropagation] For query, key, and value weights $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$ and bias $\mathbf{b}_Q, \mathbf{b}_K, \mathbf{b}_V \in \mathbb{R}^{D_{\text{aux}}}$, the backpropagation layer corresponding to a self-attention layer with H_{aux} attention heads takes a sequence $\{\partial \mathbf{y}_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and $\{\mathbf{x}_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ as input and outputs $\{\partial \mathbf{x}_t\}_{t \leq T_{\text{aux}}}$, with

$$\begin{aligned}\partial \mathbf{x}_t &= \mathbf{W}_Q^\top \partial \mathbf{q}_t + \mathbf{W}_K^\top \partial \mathbf{k}_t + \mathbf{W}_V^\top \partial \mathbf{v}_t, \quad \text{with} \\ \partial \mathbf{q}_t &= \text{VECTORIZE}(\{\sum_j a_{t,j}^h ((\partial \mathbf{y}_t^h)^\top \mathbf{v}_j^h) [\mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h]\}_{h \leq H_{\text{aux}}}); \\ \partial \mathbf{k}_t &= \text{VECTORIZE}(\{\sum_j a_{j,t}^h \mathbf{q}_j^h [(\partial \mathbf{y}_j^h)^\top (\mathbf{v}_t^h - \sum_{j'} a_{j,j'}^h \mathbf{v}_{j'}^h)]\}_{h \leq H_{\text{aux}}}); \\ \partial \mathbf{v}_t &= \text{VECTORIZE}(\{\sum_j a_{j,t}^h \partial \mathbf{y}_j^h\}_{h \leq H_{\text{aux}}})\end{aligned}$$

Here, $\mathbf{q}_t, \mathbf{k}_t$, and \mathbf{v}_t refer to query, key, and value vectors at each position t , with the attention scores $\{a_{t,j}^h\}_{t,j \leq T_{\text{aux}}, h \leq H_{\text{aux}}}$.

Derivation of gradient in Definition E.2. Recalling the definition of \mathbf{y}_t from Definition E.1,

$$\begin{aligned}\mathbf{y}_t &= \text{VECTORIZE}(\{\sum_{j \leq T_{\text{aux}}} a_{t,j}^h \mathbf{v}_j^h\}_{h \leq H_{\text{aux}}}); \quad a_{t,j}^h = \text{softmax}(\mathbf{K}^h \mathbf{q}_t^h)_j, \\ \mathbf{q}_t &= \mathbf{W}_Q \mathbf{x}_t + \mathbf{b}_Q \quad \mathbf{k}_t = \mathbf{W}_K \mathbf{x}_t + \mathbf{b}_K, \quad \mathbf{v}_t = \mathbf{W}_V \mathbf{x}_t + \mathbf{b}_V.\end{aligned}$$

$\mathbf{q}_t^h, \mathbf{k}_t^h, \mathbf{v}_t^h$ denote $\text{SPLIT}_{H_{\text{aux}}}(\mathbf{q}_t)_h$, $\text{SPLIT}_{H_{\text{aux}}}(\mathbf{k}_t)_h$, and $\text{SPLIT}_{H_{\text{aux}}}(\mathbf{v}_t)_h$ respectively for all $t \leq T_{\text{aux}}$, and $h \leq H_{\text{aux}}$. $\mathbf{K}^h \in \mathbb{R}^{T_{\text{aux}} \times D_{\text{aux}}}$ is defined with its rows as $\{\mathbf{k}_t^h\}_{t \leq T_{\text{aux}}}$ for all $h \leq H_{\text{aux}}$.

We explain the proof for an arbitrary token position t . With the application of the chain rule, we have

$$\begin{aligned}\partial \mathbf{x}_t &= \left(\frac{\partial \mathbf{q}_t}{\partial \mathbf{x}_t}\right)^\top \partial \mathbf{q}_t + \left(\frac{\partial \mathbf{k}_t}{\partial \mathbf{x}_t}\right)^\top \partial \mathbf{k}_t + \left(\frac{\partial \mathbf{v}_t}{\partial \mathbf{x}_t}\right)^\top \partial \mathbf{v}_t \\ &= \mathbf{W}_Q^\top \partial \mathbf{q}_t + \mathbf{W}_K^\top \partial \mathbf{k}_t + \mathbf{W}_V^\top \partial \mathbf{v}_t,\end{aligned}$$

where the second step follows from the definitions of $\mathbf{q}_t, \mathbf{k}_t$, and \mathbf{v}_t respectively.

Computation of $\partial \mathbf{q}_t$: With the SPLIT operation of \mathbf{q}_t across H_{aux} heads for the computation of \mathbf{y}_t , the computation of the backpropagated gradient $\partial \mathbf{q}_t$ itself needs to be split across H_{aux} heads. Furthermore, query vector \mathbf{q}_t only affects \mathbf{y}_t , implying $\frac{\partial \mathbf{y}_{t'}}{\partial \mathbf{q}_t} = 0$ for any $t' \neq t$. Thus, we have for any head $h \leq H_{\text{aux}}$, if \mathbf{y}_t^h represents the output of attention head h , given by $\sum_{j \leq T_{\text{aux}}} a_{t,j}^h \mathbf{v}_j^h$,

$$\begin{aligned}\partial \mathbf{q}_t^h &= \left(\frac{\partial \mathbf{y}_t^h}{\partial \mathbf{q}_t^h}\right)^\top \partial \mathbf{y}_t^h \\ &= \sum_{j \leq T_{\text{aux}}} \langle \mathbf{v}_j^h, \partial \mathbf{y}_t^h \rangle \frac{\partial a_{t,j}^h}{\partial \mathbf{q}_t^h} \\ &= \sum_{j \leq T_{\text{aux}}} \langle \mathbf{v}_j^h, \partial \mathbf{y}_t^h \rangle \frac{\partial}{\partial \mathbf{q}_t^h} \left(\frac{e^{\langle \mathbf{k}_j^h, \mathbf{q}_t^h \rangle}}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_t^h \rangle}} \right) \quad (8)\end{aligned}$$

$$\begin{aligned}&= \sum_{j \leq T_{\text{aux}}} \langle \mathbf{v}_j^h, \partial \mathbf{y}_t^h \rangle \left[\frac{1}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_t^h \rangle}} \frac{\partial e^{\langle \mathbf{k}_j^h, \mathbf{q}_t^h \rangle}}{\partial \mathbf{q}_t^h} - \left(\frac{e^{\langle \mathbf{k}_j^h, \mathbf{q}_t^h \rangle}}{(\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_t^h \rangle})^2} \right) \sum_{j' \leq T_{\text{aux}}} \frac{\partial e^{\langle \mathbf{k}_{j'}^h, \mathbf{q}_t^h \rangle}}{\partial \mathbf{q}_t^h} \right] \quad (9)\end{aligned}$$

$$\begin{aligned}&= \sum_{j \leq T_{\text{aux}}} \langle \mathbf{v}_j^h, \partial \mathbf{y}_t^h \rangle \left[\left(\frac{e^{\langle \mathbf{k}_j^h, \mathbf{q}_t^h \rangle}}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_t^h \rangle}} \right) \mathbf{k}_j^h - \left(\frac{e^{\langle \mathbf{k}_j^h, \mathbf{q}_t^h \rangle}}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_t^h \rangle}} \right) \sum_{j' \leq T_{\text{aux}}} \left(\frac{e^{\langle \mathbf{k}_{j'}^h, \mathbf{q}_t^h \rangle}}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_t^h \rangle}} \right) \mathbf{k}_{j'}^h \right] \quad (10)\end{aligned}$$

$$= \sum_{j \leq T_{\text{aux}}} a_{t,j}^h \langle \mathbf{v}_j^h, \partial_{\mathbf{y}_t^h} \rangle \left(\mathbf{k}_j^h - \sum_{j' \leq T_{\text{aux}}} a_{t,j'}^h \mathbf{k}_{j'}^h \right).$$

In Equation (8), we have expanded the definition of softmax in $a_{t,j}^h := \text{softmax}(\mathbf{K}^h \mathbf{q}_t^h)_j$ in order to better motivate the derivative of $a_{t,j}^h$ w.r.t. \mathbf{q}_t^h . Finally, $\partial_{\mathbf{q}_t}$ is given by $\text{VECTORIZE}(\{\partial_{\mathbf{q}_t^h}\}_{h \leq H_{\text{aux}}})$.

Computation of $\partial_{\mathbf{k}_t}$: Continuing as the computation of $\partial_{\mathbf{q}_t}$, we split the computation of $\partial_{\mathbf{k}_t}$ across the H_{aux} attention heads. However, unlike \mathbf{q}_t , \mathbf{k}_t affects \mathbf{y}_j for all $j \leq T_{\text{aux}}$. For any head $h \leq H_{\text{aux}}$, we follow the chain-rule step by step to get

$$\begin{aligned} \partial_{\mathbf{k}_t^h} &= \sum_{j \leq T_{\text{aux}}} \left(\frac{\partial \mathbf{y}_j^h}{\partial \mathbf{k}_t^h} \right)^\top \partial_{\mathbf{y}_j^h} = \sum_{j \leq T_{\text{aux}}} \left(\frac{\partial \sum_{j' \leq T_{\text{aux}}} a_{j,j'}^h \mathbf{v}_{j'}^h}{\partial \mathbf{k}_t^h} \right)^\top \partial_{\mathbf{y}_j^h} \\ &= \sum_{j \leq T_{\text{aux}}} \langle \mathbf{v}_t^h, \partial_{\mathbf{y}_j^h} \rangle \frac{\partial a_{j,t}^h}{\partial \mathbf{k}_t^h} + \sum_{j \leq T_{\text{aux}}} \sum_{j' \leq T_{\text{aux}}; j' \neq t} \langle \mathbf{v}_{j'}^h, \partial_{\mathbf{y}_j^h} \rangle \frac{\partial a_{j,j'}^h}{\partial \mathbf{k}_t^h} \end{aligned} \quad (11)$$

$$= \sum_{j \leq T_{\text{aux}}} \langle \mathbf{v}_t^h, \partial_{\mathbf{y}_j^h} \rangle \frac{\partial}{\partial \mathbf{k}_t^h} \left(\frac{e^{\langle \mathbf{k}_t^h, \mathbf{q}_j^h \rangle}}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_j^h \rangle}} \right) \quad (12)$$

$$+ \sum_{j \leq T_{\text{aux}}} \sum_{j' \leq T_{\text{aux}}; j' \neq t} \langle \mathbf{v}_{j'}^h, \partial_{\mathbf{y}_j^h} \rangle \frac{\partial}{\partial \mathbf{k}_t^h} \left(\frac{e^{\langle \mathbf{k}_{j'}^h, \mathbf{q}_j^h \rangle}}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_j^h \rangle}} \right) \quad (13)$$

$$= \sum_{j \leq T_{\text{aux}}} \langle \mathbf{v}_t^h, \partial_{\mathbf{y}_j^h} \rangle \left[\left(\frac{e^{\langle \mathbf{k}_t^h, \mathbf{q}_j^h \rangle}}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_j^h \rangle}} \right) \mathbf{q}_j^h - \left(\frac{e^{\langle \mathbf{k}_t^h, \mathbf{q}_j^h \rangle}}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_j^h \rangle}} \right)^2 \mathbf{q}_j^h \right] \quad (14)$$

$$- \sum_{j \leq T_{\text{aux}}} \sum_{j' \leq T_{\text{aux}}; j' \neq t} \langle \mathbf{v}_{j'}^h, \partial_{\mathbf{y}_j^h} \rangle \left(\frac{e^{\langle \mathbf{k}_{j'}^h, \mathbf{q}_j^h \rangle}}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_j^h \rangle}} \right) \left(\frac{e^{\langle \mathbf{k}_t^h, \mathbf{q}_j^h \rangle}}{\sum_{t' \leq T_{\text{aux}}} e^{\langle \mathbf{k}_{t'}^h, \mathbf{q}_j^h \rangle}} \right) \mathbf{q}_j^h \quad (15)$$

$$\begin{aligned} &= \sum_{j \leq T_{\text{aux}}} \langle \mathbf{v}_t^h, \partial_{\mathbf{y}_j^h} \rangle (a_{j,t}^h - (a_{j,t}^h)^2) \mathbf{q}_j^h - \sum_{j \leq T_{\text{aux}}} \sum_{j' \leq T_{\text{aux}}; j' \neq t} \langle \mathbf{v}_{j'}^h, \partial_{\mathbf{y}_j^h} \rangle a_{j,j'}^h a_{j,t}^h \mathbf{q}_j^h \\ &= \sum_{j \leq T_{\text{aux}}} a_{j,t}^h \langle \partial_{\mathbf{y}_j^h} \mathbf{v}_t^h, \mathbf{v}_t^h \rangle - \sum_{j'} a_{j,j'}^h \langle \mathbf{v}_{j'}^h, \mathbf{v}_t^h \rangle \mathbf{q}_j^h \end{aligned}$$

In Equation (11), we separate the inside sum into two components, since the derivative w.r.t. \mathbf{k}_t^h differ for the two components, as outlined in the derivation of Equation (14) from Equation (12), and Equation (15) from Equation (13). We have skipped a step going from Equations (12) and (13) to Equations (14) and (15) due to typographical simplicity. The skipped step is extremely similar to Equation (9) in the derivation of $\partial_{\mathbf{q}_t^h}$. Finally, $\partial_{\mathbf{k}_t}$ is given by $\text{VECTORIZE}(\{\partial_{\mathbf{k}_t^h}\}_{h \leq H_{\text{aux}}})$.

Computation of $\partial_{\mathbf{v}_t}$: Similar to the gradient computation of \mathbf{q}_t , the computation of $\partial_{\mathbf{v}_t}$ needs to be split across the H_{aux} attention heads. However, like \mathbf{k}_t , \mathbf{v}_t affects \mathbf{y}_j for all $j \leq T_{\text{aux}}$. For any head $h \leq H_{\text{aux}}$, we follow the chain-rule step by step to get

$$\partial_{\mathbf{v}_t^h} = \sum_{j \leq T_{\text{aux}}} \left(\frac{\partial \mathbf{y}_j^h}{\partial \mathbf{v}_t^h} \right)^\top \partial_{\mathbf{y}_j^h} = \sum_{j \leq T_{\text{aux}}} \left(\frac{\partial \sum_{j' \leq T_{\text{aux}}} a_{j,j'}^h \mathbf{v}_{j'}^h}{\partial \mathbf{v}_t^h} \right)^\top \partial_{\mathbf{y}_j^h} = \sum_{j \leq T_{\text{aux}}} a_{j,t}^h \partial_{\mathbf{y}_j^h}$$

□

Theorem E.4. With the notations in Definitions A.7, E.1 and E.2, if on a given input sequence $\{\mathbf{x}_t\}_{t=1}^{T_{\text{aux}}}$, with its query, key, and value vectors $\{\mathbf{q}_t, \mathbf{k}_t, \mathbf{v}_t\}_{t=1}^{T_{\text{aux}}}$, all the H_{aux} attention heads are ε -hard for some $\varepsilon > 0$, then for a given sequence of gradients $\{\partial_{\mathbf{y}_t}\}_{t=1}^{T_{\text{aux}}}$,

$$\|\partial_{\mathbf{q}_t}\|_2, \|\partial_{\mathbf{k}_t}\|_2 \leq \mathcal{O}(\varepsilon B_x^2 B_w^2 B_y), \quad \text{for all } t \leq T_{\text{aux}},$$

where $B_x = \max_{t \leq T_{\text{aux}}} \|\mathbf{x}_t\|_2$, $B_y = \max_{t \leq T_{\text{aux}}} \|\partial \mathbf{y}_t\|_2$, and $B_w = \max\{\|\mathbf{W}_K\|_2, \|\mathbf{W}_Q\|_2, \|\mathbf{W}_V\|_2, \|\mathbf{b}_V\|_2, \|\mathbf{b}_K\|_2, \|\mathbf{b}_V\|_2\}$.

This implies, for each position t , $\|\widehat{\partial \mathbf{x}_t} - \partial \mathbf{x}_t\|_2 \leq \mathcal{O}(\varepsilon B_x^2 B_w^3 B_y)$.

Proof of Theorem E.4. For typographical simplicity, we discuss the proof at an arbitrary position t . Recall the definition of an ε -hard attention head from Definition E.3. An attention head is defined to be ε -hard on an input sequence $\{\mathbf{x}_t\}_{t=1}^{T_{\text{aux}}}$, if for each position t , there exists a position t_0 such that the attention score $a_{t,t_0} \geq 1 - \varepsilon$.

For the proof, we simply focus on $\partial_{\mathbf{q}_t}$, and the proof for $\partial_{\mathbf{k}_t}$ follows like-wise.

Bounds on \mathbf{q}_t : Recalling the definition of $\partial_{\mathbf{q}_t}$ from Definition E.2, we have

$$\partial_{\mathbf{q}_t} = \text{VECTORIZE}(\{\sum_j a_{t,j}^h ((\partial \mathbf{y}_t^h)^\top \mathbf{v}_j^h) [\mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h]\}_{h \leq H_{\text{aux}}}).$$

Focusing on a head $h \leq H_{\text{aux}}$, define $\partial_{\mathbf{q}_t^h} = \sum_j a_{t,j}^h ((\partial \mathbf{y}_t^h)^\top \mathbf{v}_j^h) [\mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h]$ and $t_0 \leq T_{\text{aux}}$ as the token position where the \mathbf{q}_t attends the most to, i.e. $a_{t,t_0} \geq 1 - \varepsilon$ and $\sum_{j \leq T_{\text{aux}}; j \neq t_0} a_{t,j}^h \leq \varepsilon$. Then,

$$\begin{aligned} \|\partial_{\mathbf{q}_t^h}\|_2 &= \left\| \sum_j a_{t,j}^h ((\partial \mathbf{y}_t^h)^\top \mathbf{v}_j^h) [\mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h] \right\|_2 \\ &= \left\| a_{t,t_0}^h ((\partial \mathbf{y}_t^h)^\top \mathbf{v}_{t_0}^h) [\mathbf{k}_{t_0}^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h] + \sum_{j \neq t_0} a_{t,j}^h ((\partial \mathbf{y}_t^h)^\top \mathbf{v}_j^h) [\mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h] \right\|_2 \\ &\leq \underbrace{\left\| a_{t,t_0}^h ((\partial \mathbf{y}_t^h)^\top \mathbf{v}_{t_0}^h) [\mathbf{k}_{t_0}^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h] \right\|_2}_{\text{Term1}} + \underbrace{\left\| \sum_{j \neq t_0} a_{t,j}^h ((\partial \mathbf{y}_t^h)^\top \mathbf{v}_j^h) [\mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h] \right\|_2}_{\text{Term2}}, \end{aligned}$$

where the final step uses a Cauchy-Schwartz inequality. We focus on the two terms separately.

1. Term1: Focusing on $\mathbf{k}_{t_0}^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h$, we have

$$\begin{aligned} \left\| \mathbf{k}_{t_0}^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h \right\|_2 &= \left\| (1 - a_{t,t_0}) \mathbf{k}_{t_0}^h - \sum_{j' \neq t_0} a_{t,j'}^h \mathbf{k}_{j'}^h \right\|_2 \\ &\leq (1 - a_{t,t_0}) \|\mathbf{k}_{t_0}^h\|_2 + \sum_{j' \neq t_0} a_{t,j'}^h \|\mathbf{k}_{j'}^h\|_2 \\ &\leq ((1 - a_{t,t_0}) + \sum_{j' \neq t_0} a_{t,j'}^h) \max_j \|\mathbf{k}_j^h\|_2 \\ &\leq 2\varepsilon \max_j \|\mathbf{k}_j^h\|_2. \end{aligned} \tag{16}$$

We use a Cauchy-Schwartz inequality in the second and third steps and the attention head behavior in the final step.

Hence, Term1 can now be bounded as follows:

$$\begin{aligned} \left\| a_{t,t_0}^h ((\partial \mathbf{y}_t^h)^\top \mathbf{v}_{t_0}^h) [\mathbf{k}_{t_0}^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h] \right\|_2 &= a_{t,t_0}^h \left| (\partial \mathbf{y}_t^h)^\top \mathbf{v}_{t_0}^h \right| \left\| \mathbf{k}_{t_0}^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h \right\|_2 \\ &\leq 2\varepsilon \left\| \partial \mathbf{y}_t^h \right\|_2 \left\| \mathbf{v}_{t_0}^h \right\|_2 \max_j \|\mathbf{k}_j^h\|_2. \end{aligned}$$

In the final step, in addition to the bound from Equation (16), we use a Cauchy-Schwartz inequality to bound $\left| (\partial \mathbf{y}_t^h)^\top \mathbf{v}_{t_0}^h \right|$ and bound the attention score a_{t,t_0}^h by 1.

2. Term2: Focusing on $\mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h$ for any $j \leq T_{\text{aux}}$, we have using two Cauchy-Schwartz inequalities:

$$\left\| \mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h \right\|_2 \leq \|\mathbf{k}_j^h\|_2 + \left\| \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h \right\|_2 \leq (1 + \sum_{j'} a_{t,j'}^h) \max_{j'} \|\mathbf{k}_{j'}^h\|_2 = 2 \max_{j'} \|\mathbf{k}_{j'}^h\|_2. \quad (17)$$

Hence,

$$\begin{aligned} \left\| \sum_{j \neq t_0} a_{t,j}^h ((\partial_{\mathbf{y}_t^h})^\top \mathbf{v}_j^h) [\mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h] \right\|_2 &\leq \left(\sum_{j \neq t_0} a_{t,j}^h \right) \max_j |(\partial_{\mathbf{y}_t^h})^\top \mathbf{v}_j^h| \left\| \mathbf{k}_j^h - \sum_{j'} a_{t,j'}^h \mathbf{k}_{j'}^h \right\|_2 \\ &\leq 2\varepsilon \|\partial_{\mathbf{y}_t^h}\|_2 \left(\max_j \|\mathbf{v}_j^h\|_2 \right) \left(\max_{j'} \|\mathbf{k}_{j'}^h\|_2 \right). \end{aligned}$$

In the final step, in addition to the bound from Equation (17), we use a Cauchy-Schwartz inequality to bound $|(\partial_{\mathbf{y}_t^h})^\top \mathbf{v}_j^h|$ and use the ε -hard behavior of the attention head to bound $\sum_{j \neq t_0} a_{t,j}^h$.

Combining the bounds on both terms, we have

$$\begin{aligned} \|\partial_{\mathbf{q}_t^h}\|_2 &\leq 2\varepsilon \|\partial_{\mathbf{y}_t^h}\|_2 \|\mathbf{v}_{t_0}^h\|_2 \max_j \|\mathbf{k}_j^h\|_2 + 2\varepsilon \|\partial_{\mathbf{y}_t^h}\|_2 \left(\max_j \|\mathbf{v}_j^h\|_2 \right) \left(\max_{j'} \|\mathbf{k}_{j'}^h\|_2 \right) \\ &\leq 4\varepsilon \|\partial_{\mathbf{y}_t^h}\|_2 \left(\max_j \|\mathbf{v}_j^h\|_2 \right) \left(\max_{j'} \|\mathbf{k}_{j'}^h\|_2 \right). \end{aligned}$$

We bound the remaining terms as follows.

- $\|\partial_{\mathbf{y}_t^h}\|_2 \leq B_y$, under the bounded assumption of the gradients.
- For any $j \leq T_{\text{aux}}$, we have $\|\mathbf{k}_j^h\|_2 \leq \|\mathbf{k}_j\|_2$ since $\mathbf{k}_j = \text{VECTORIZE}(\{\mathbf{k}_j^{h'}\}_{h' \in H_{\text{aux}}})$. Furthermore, from the definition of the key vector \mathbf{k}_j , $\|\mathbf{k}_j\|_2 = \|\mathbf{W}_K \mathbf{x}_j + \mathbf{b}_K\|_2 \leq \|\mathbf{W}_K\|_2 \|\mathbf{x}_j\|_2 + \|\mathbf{b}_K\|_2$ with a Cauchy-Schwartz inequality. Under the bounded assumptions of \mathbf{W}_K , \mathbf{b}_K and input \mathbf{x}_j , we have $\|\mathbf{k}_j\|_2 \leq B_w(1 + B_x)$.
- Similar procedure can be followed for bounding $\max_j \|\mathbf{v}_j^h\|_2$.

Thus, we have $\|\partial_{\mathbf{q}_t^h}\|_2 \leq 4\varepsilon \|\partial_{\mathbf{y}_t^h}\|_2 \left(\max_j \|\mathbf{v}_j^h\|_2 \right) \left(\max_{j'} \|\mathbf{k}_{j'}^h\|_2 \right) \leq 4\varepsilon B_w^2 (1 + B_x)^2 B_y$.

Bounds on $\|\widehat{\partial_{\mathbf{x}_t}} - \partial_{\mathbf{x}_t}\|_2$: From the definitions of $\widehat{\partial_{\mathbf{x}_t}}$ and $\partial_{\mathbf{x}_t}$ from Definitions A.7 and E.6, we have

$$\begin{aligned} \|\widehat{\partial_{\mathbf{x}_t}} - \partial_{\mathbf{x}_t}\|_2 &= \|\mathbf{W}_K^\top \partial_{\mathbf{k}_t} + \mathbf{W}_Q^\top \partial_{\mathbf{q}_t}\|_2 \leq \|\mathbf{W}_K\|_2 \|\partial_{\mathbf{k}_t}\|_2 + \|\mathbf{W}_Q\|_2 \|\partial_{\mathbf{q}_t}\|_2 \\ &\leq 8\varepsilon B_w^3 (1 + B_x)^2 B_y = \mathcal{O}(\varepsilon B_w^3 B_x^2 B_y), \end{aligned}$$

where we use Cauchy-schwartz inequality in the second step. We use the assumed bounds on $\|\mathbf{W}_Q\|_2$, $\|\mathbf{W}_K\|_2$, and the computed bounds on $\|\partial_{\mathbf{q}_t}\|_2$, $\|\partial_{\mathbf{k}_t}\|_2$ in the pre-final step. \square

F Layer normalization

We repeat the definition of layer normalization from the main paper below.

Definition A.1. [Layer Normalization] Define a normalization function $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that performs $f(\mathbf{x}) = (\mathbf{x} - \mu)/\sigma$, where μ and σ are the mean and standard deviation of \mathbf{x} , respectively. Then, layer normalization with parameters $\gamma, \mathbf{b} \in \mathbb{R}^{D_{\text{aux}}}$ takes as input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ and outputs $\mathbf{y} \in \mathbb{R}^{D_{\text{aux}}}$, which is computed as $\mathbf{z} = f(\mathbf{x}), \mathbf{y} = \gamma \odot \mathbf{z} + \mathbf{b}$.

In the discussions below, we consider a layer normalization layer in the auxiliary model with parameters $\{\gamma, \mathbf{b}\}$ that takes in input sequence $\mathbf{x}_1, \dots, \mathbf{x}_{T_{\text{aux}}}$ and outputs $\mathbf{y}_1, \dots, \mathbf{y}_{T_{\text{aux}}}$, with $\mathbf{y}_t = \gamma \odot \mathbf{z}_t + \mathbf{b}; \mathbf{z}_t = f(\mathbf{x}_t)$ for each $t \leq T_{\text{aux}}$. Since this involves a token-wise operation, we will present our constructed modules with a general token position t and the prefix tokens $\{\mathbf{v}_j\}$. We will use \mathbf{W}_γ as a diagonal matrix in $\mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$, containing γ on its main diagonal.

TINT Layer normalization Forward module The input embedding to this module \mathbf{e}_t will contain \mathbf{x}_t in its first D_{aux} coordinates. The layer normalization computation can be divided into two sub-operations: (a) application of f , and (b) linear computation using γ, \mathbf{b} . We will present a TINT module for each sub-operation.

We can represent the function f using a layer normalization operation itself, with its weight and bias parameters set as $\mathbf{1}$ and $\mathbf{0}$ respectively. However, since the relevant input exists only in the first D_{aux} coordinates, the operation on the first D_{aux} coordinates needs to be independent of the rest of the coordinates. To do so, we instead use Group normalization (Definition F.3) on \mathbf{e}_t , with groups of size D_{aux} .

Now, the embedding \mathbf{e}_t contains $f(\mathbf{x}_t)$ in its first D_{aux} coordinates. The second sub-operation can then be viewed as a Linear Layer computation, i.e. $\mathbf{y}_t = \mathbf{W}_\gamma \mathbf{x}_t + \mathbf{b}$. Hence, we simply stack the rows of \mathbf{W}_γ and \mathbf{b}_γ onto the prefix tokens $\{\mathbf{v}_j\}$ and call the TINT Linear Forward module (Appendix D).

Auxiliary’s gradient backpropagation through layer normalization With the definition of layer normalization and the normalization function f in Definition A.1, the auxiliary’s backpropagation operation takes in the loss gradient w.r.t. output ($\partial_{\mathbf{y}}$) and computes the loss gradient w.r.t. input ($\partial_{\mathbf{x}}$).

Definition A.2. [Exact Gradient for Layer Normalization] Using notations in Definition A.1, given the gradient of the loss w.r.t the output of the Layer Normalization $\partial_{\mathbf{y}}$, backpropagation computes $\partial_{\mathbf{x}}$ as

$$\partial_{\mathbf{x}} = (\partial_{\mathbf{z}} - D_{\text{aux}}^{-1} \sum_{i=1}^{D_{\text{aux}}} \partial_{z_i} - \langle \partial_{\mathbf{z}}, \mathbf{z} \rangle \mathbf{z}) / \sigma \quad \partial_{\mathbf{z}} = \gamma \odot \partial_{\mathbf{y}}.$$

Complexity of true backpropagation The above operation is computation heavy since it involves computing (a) $\partial_{\mathbf{z}}$, (b) $f(\partial_{\mathbf{z}})$, (c) $\langle \partial_{\mathbf{z}}, \mathbf{z} \rangle \mathbf{z}$, and (d) multiplying by a factor of $\frac{1}{\sigma}$. $\langle \partial_{\mathbf{z}}, \mathbf{z} \rangle \mathbf{z}$ in itself will require two MLP layers, following Lemma C.4. In order to reduce the number of layers, we turn to first-order Taylor expansion for approximating the above operation.

Definition A.3. [ϵ -approximate Layer Normalization Gradient] With notations defined above, this layer takes $\partial_{\mathbf{y}}, \mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ as input and outputs $\widehat{\partial}_{\mathbf{x}} = \frac{1}{\epsilon} (f(\mathbf{x} + \epsilon \gamma \odot \partial_{\mathbf{y}}) - f(\mathbf{x}))$.

The following theorem shows that the first-order gradient is a good approximation of the true gradient, and in the limit of ϵ tending to 0, the approximation error tends to 0 as well.

Theorem F.1. For any $\epsilon > 0$, and a layer normalization layer with parameters $\gamma, \mathbf{b} \in \mathbb{R}^{D_{\text{aux}}}$, for an input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ and gradient $\partial_{\mathbf{y}} \in \mathbb{R}^{D_{\text{aux}}}$,

$$\left\| \widehat{\partial}_{\mathbf{x}} - \partial_{\mathbf{x}} \right\|_2 \leq \mathcal{O}(\epsilon D_{\text{aux}}^{3/2} \sigma^{-2} \|\gamma\|_2^2 \|\partial_{\mathbf{y}}\|_2^2),$$

where σ denotes the standard deviation of \mathbf{x} . $\partial_{\mathbf{x}}, \widehat{\partial}_{\mathbf{x}}$ have been computed from $\mathbf{x}, \partial_{\mathbf{y}}$ and ϵ using Definitions A.2 and A.3.

TINT Layer normalization backpropagation module The input embeddings \mathbf{e}_t contain $\partial_{\mathbf{y}_t}$ at each position t in the first D_{aux} coordinates. Since we further need the input to the auxiliary’s layer normalization layer under consideration, we copy \mathbf{x}_t from the TINT Layer normalization Forward module at each position t using residual connections. Furthermore, residual connections have been used to copy the contents of the prefix tokens $\{\mathbf{v}_j\}$ from the Layer normalization Forward module, which contain $\mathbf{W}_\gamma, \mathbf{b}$. Recall that for ease of presentation, we use \mathbf{z}_t to represent $f(\mathbf{x}_t)$.

We set ϵ as a hyperparameter and return $\widehat{\partial_x}$ as the output of this module. The computation of $\widehat{\partial_x}$ can be divided into two sub-operations: (a) computation of $\partial_{z_t} := \gamma \odot \partial_{y_t}$, and (b) computation of $\frac{1}{\epsilon}(f(\mathbf{x}_t + \epsilon \partial_{z_t}) - f(\mathbf{x}_t))$. We represent each sub-operation as a TINT module.

To compute $\partial_{z_t} := \gamma \odot \partial_{y_t} = W_\gamma \partial_{y_t}$, we can observe that the required operation is identical to backpropagating through a linear layer with parameters W_γ and \mathbf{b} . Hence, we simply call the Linear Backpropagation module on the current embeddings. We use residual connections to retain \mathbf{x}_t at each location t , and the contents of the prefix tokens $\{v_j\}$.

Now, the embedding e_t contains ∂_{z_t} and \mathbf{x}_t . In order to backpropagate through f , we first use a linear layer to compute $\mathbf{x}_t + \epsilon \partial_{z_t}$ and retain \mathbf{x}_t . Following the same procedure as the Forward module, we use a Group normalization layer with weight and bias parameters $\mathbf{1}$ and $\mathbf{0}$ respectively, to compute $f(\mathbf{x}_t + \epsilon \partial_{z_t})$ and $f(\mathbf{x}_t)$. Finally, we use a linear layer to compute $\frac{1}{\epsilon}(f(\mathbf{x}_t + \epsilon \partial_{z_t}) - f(\mathbf{x}_t))$.

Auxiliary’s Descent update And finally, the auxiliary’s descent operation updates parameters γ, \mathbf{b} using a batch of inputs $\{\mathbf{x}_t\}_{t \leq T}$ and the loss gradient w.r.t. the corresponding outputs $\{\partial_{y_t}\}_{t \leq T}$.

Definition F.2 (Auxiliary’s layer normalization descent). For parameters $\gamma, \mathbf{b} \in \mathbb{R}^{D_{\text{aux}}}$, descent update takes in a batch of inputs $\{\mathbf{x}_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and gradients $\{\partial_{y_t} \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and updates the parameters as follows:

$$\gamma \leftarrow \gamma - \eta \sum_{t \leq T_{\text{aux}}} \partial_{y_t} \odot \mathbf{z}_t; \quad \mathbf{b} \leftarrow \mathbf{b} - \eta \sum_{t \leq T_{\text{aux}}} \partial_{y_t},$$

where \mathbf{z}_t represents $f(\mathbf{x}_t)$.

The update of γ involves an elementwise multiplication between ∂_{y_t} and \mathbf{z}_t , which requires an MLP layer (Lemma C.4). With the prefix tokens containing the rows of W_γ and \mathbf{b} , we instead consider the update of \mathbf{b} alone with the descent update.

TINT Layer normalization descent module The input embeddings contain ∂_{y_t} in the first D_{aux} coordinates. The prefix tokens contain W_γ, \mathbf{b} , which have been copied from the Forward module using residual connections. The update of \mathbf{b} is identical to the auxiliary’s descent update through a linear layer. Hence, we apply a TINT Linear descent module to the current embeddings, updating only the bias \mathbf{b} and switching off the update to W_γ .

F.1 Additional definitions

We describe TINT group normalization layer below, which we use in different modules to simulate the auxiliary’s layer normalization operations.

Definition F.3 (TINT D_{aux} -Group normalization). Define a normalization function $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that performs $f(\mathbf{x}) = (\mathbf{x} - \mu)/\sigma$, where μ and σ are the mean and standard deviation of \mathbf{x} , respectively. Then, D_{aux} -Group RMSnorm with parameters $\gamma^{\text{TINT}}, \mathbf{b}^{\text{TINT}} \in \mathbb{R}^{D_{\text{aux}}}$ takes as input $\mathbf{x} \in \mathbb{R}^{D_{\text{sim}}}$ and outputs $\mathbf{y} = \text{VECTORIZE}(\{\mathbf{y}^h \in \mathbb{R}^{D_{\text{aux}}}\}_{h \leq \lfloor D_{\text{sim}}/D_{\text{aux}} \rfloor})$, with

$$\mathbf{y}^h = \gamma^{\text{TINT}} \odot f(\mathbf{x}^h) + \mathbf{b}^{\text{TINT}},$$

where $\mathbf{x}^h = \text{SPLIT}_{\lfloor D_{\text{sim}}/D_{\text{aux}} \rfloor}(\mathbf{x})_h$.

F.2 Proof of theorems and gradient definitions

We restate the theorems and definitions, before presenting their proofs for easy referencing.

Definition A.2. [Exact Gradient for Layer Normalization] Using notations in Definition A.1, given the gradient of the loss w.r.t the output of the Layer Normalization $\partial_{\mathbf{y}}$, backpropagation computes $\partial_{\mathbf{x}}$ as

$$\partial_{\mathbf{x}} = (\partial_{\mathbf{z}} - D_{\text{aux}}^{-1} \sum_{i=1}^{D_{\text{aux}}} \partial_{z_i} - \langle \partial_{\mathbf{z}}, \mathbf{z} \rangle \mathbf{z}) / \sigma \quad \partial_{\mathbf{z}} = \gamma \odot \partial_{\mathbf{y}}.$$

Derivation of gradient in Definition A.2. With the normalization function f and parameters $\mathbf{x}, \mathbf{b} \in \mathbb{R}^{D_{\text{aux}}}$, recall from Definition A.1 that given an input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$, a layer normalization layer returns

$\mathbf{y} = \gamma \odot \mathbf{z} + \mathbf{b}$; $\mathbf{z} = f(\mathbf{x})$. Let μ and σ denote the mean and standard deviation of \mathbf{x} . They can be computed as

$$\mu = \frac{1}{D_{\text{aux}}} \sum_{i=1}^{D_{\text{aux}}} x_i, \quad \sigma = \sqrt{\frac{1}{D_{\text{aux}}} \sum_{i=1}^{D_{\text{aux}}} (x_i - \mu)^2}.$$

With the chain rule, we can compute $\partial_{\mathbf{x}}$ from $\partial_{\mathbf{y}}$ as follows.

$$\partial_{\mathbf{x}} = \left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right)^\top \partial_{\mathbf{z}}; \quad \text{with } \partial_{\mathbf{z}} = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{z}} \right)^\top \partial_{\mathbf{y}}. \quad (18)$$

Since $\mathbf{y} = \gamma \odot \mathbf{z} + \mathbf{b}$, we have $\frac{\partial \mathbf{y}}{\partial \mathbf{z}} = \mathbf{W}_\gamma$, where \mathbf{W}_γ represents a diagonal matrix with γ on the main diagonal. Thus, $\partial_{\mathbf{z}} = \mathbf{W}_\gamma \partial_{\mathbf{y}} = \gamma \odot \partial_{\mathbf{y}}$.

With $\mathbf{z} = f(\mathbf{x}) = \frac{\mathbf{x} - \mu}{\sigma}$, we have

$$\begin{aligned} \frac{\partial \mathbf{z}}{\partial \mathbf{x}} &= \frac{\partial}{\partial \mathbf{x}} \left(\frac{\mathbf{x} - \mu}{\sigma} \right) = \frac{1}{\sigma} \frac{\partial \mathbf{x}}{\partial \mathbf{x}} - \frac{1}{\sigma} \frac{\partial \mu}{\partial \mathbf{x}} - \frac{(\mathbf{x} - \mu)}{\sigma^2} \left(\frac{\partial \sigma}{\partial \mathbf{x}} \right)^\top \\ &= \frac{1}{\sigma} \left(\mathbf{I} - \frac{1}{D_{\text{aux}}} \mathbf{1} \mathbf{1}^\top - \mathbf{z} \mathbf{z}^\top \right). \end{aligned} \quad (19)$$

In the final step, we require $\frac{\partial \mu}{\partial \mathbf{x}}$ and $\frac{\partial \sigma}{\partial \mathbf{x}}$, which are computed as follows.

- $\frac{\partial \mu}{\partial \mathbf{x}} \in \mathbb{R}^{D_{\text{aux}}}$ with its j th element given by

$$\left(\frac{\partial \mu}{\partial \mathbf{x}} \right)_j = \frac{\partial \mu}{\partial x_j} = \frac{\partial}{\partial x_j} \left(\frac{1}{D_{\text{aux}}} \sum_{i=1}^{D_{\text{aux}}} x_i \right) = \frac{1}{D_{\text{aux}}}.$$

- $\frac{\partial \sigma}{\partial \mathbf{x}} \in \mathbb{R}^{D_{\text{aux}}}$ with its j th element given by

$$\begin{aligned} \left(\frac{\partial \sigma}{\partial \mathbf{x}} \right)_j &= \frac{\partial \sigma}{\partial x_j} = \frac{\partial}{\partial x_j} \left(\sqrt{\frac{1}{D_{\text{aux}}} \sum_{i=1}^{D_{\text{aux}}} (x_i - \mu)^2} \right) \\ &= \frac{1}{\sqrt{\sum_{i=1}^{D_{\text{aux}}} (x_i - \mu)^2}} \sum_{i=1}^{D_{\text{aux}}} (x_i - \mu) \frac{\partial (x_i - \mu)}{\partial x_j} \\ &= \frac{1}{\sqrt{\sum_{i=1}^{D_{\text{aux}}} (x_i - \mu)^2}} \left((x_j - \mu) - \frac{1}{D_{\text{aux}}} \sum_{i=1}^{D_{\text{aux}}} (x_i - \mu) \right) = \frac{x_j - \mu}{\sigma} := z_j, \end{aligned}$$

where we have re-utilized the $\frac{\partial \mu}{\partial \mathbf{x}}$ in the pre-final step.

Hence, from Equation (18),

$$\partial_{\mathbf{x}} = \left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right)^\top \partial_{\mathbf{z}} = \frac{1}{\sigma} \left(\mathbf{I} - \frac{1}{D_{\text{aux}}} \mathbf{1} \mathbf{1}^\top - \mathbf{z} \mathbf{z}^\top \right) \partial_{\mathbf{z}} = \frac{1}{\sigma} \left(\partial_{\mathbf{z}} - \frac{1}{D_{\text{aux}}} \langle \mathbf{1}, \partial_{\mathbf{z}} \rangle \mathbf{1} - \langle \mathbf{z}, \partial_{\mathbf{z}} \rangle \mathbf{z} \right).$$

□

We repeat Theorem F.1 for easier reference.

Theorem F.1. For any $\epsilon > 0$, and a layer normalization layer with parameters $\gamma, \mathbf{b} \in \mathbb{R}^{D_{\text{aux}}}$, for an input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ and gradient $\partial_{\mathbf{y}} \in \mathbb{R}^{D_{\text{aux}}}$,

$$\left\| \widehat{\partial_{\mathbf{x}}} - \partial_{\mathbf{x}} \right\|_2 \leq \mathcal{O}(\epsilon D_{\text{aux}}^{3/2} \sigma^{-2} \|\gamma\|_2^2 \|\partial_{\mathbf{y}}\|_2^2),$$

where σ denotes the standard deviation of \mathbf{x} . $\partial_{\mathbf{x}}, \widehat{\partial_{\mathbf{x}}}$ have been computed from $\mathbf{x}, \partial_{\mathbf{y}}$ and ϵ using Definitions A.2 and A.3.

Proof of Theorem F.1. With the normalization function f and parameters $\mathbf{x}, \mathbf{b} \in \mathbb{R}^{D_{\text{aux}}}$, recall from Definition A.1 that given an input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$, a layer normalization layer returns $\mathbf{y} = \gamma \odot \mathbf{z} + \mathbf{b}$; $\mathbf{z} = f(\mathbf{x})$. Let μ and σ denote the mean and standard deviation of \mathbf{x} . They can be computed as

$$\mu = \frac{1}{D_{\text{aux}}} \sum_{i=1}^{D_{\text{aux}}} x_i, \quad \sigma = \sqrt{\frac{1}{D_{\text{aux}}} \sum_{i=1}^{D_{\text{aux}}} (x_i - \mu)^2}.$$

We will refer to $\frac{\partial \mathbf{z}}{\partial \mathbf{x}}$ from Equation (19) and the formulation of $\partial_{\mathbf{x}}$ from Equation (18) for our current proof. To recall, they are

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \frac{1}{\sigma} \left(\mathbf{I} - \frac{1}{D_{\text{aux}}} \mathbf{1}\mathbf{1}^\top - \mathbf{z}\mathbf{z}^\top \right), \quad \partial_{\mathbf{x}} = \left(\frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right)^\top \partial_{\mathbf{z}}.$$

Using a second-order Taylor expansion of the normalization function f around \mathbf{x} , we have

$$\begin{aligned} f(\mathbf{x} + \epsilon \partial_{\mathbf{z}}) &= f(\mathbf{x}) + \epsilon \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \partial_{\mathbf{z}} + \int_0^\epsilon \partial_{\mathbf{z}}^\top \frac{\partial}{\partial \mathbf{x}_\theta} \left(\frac{\partial f(\mathbf{x}_\theta)}{\partial \mathbf{x}_\theta} \right) \partial_{\mathbf{z}} \theta d\theta \\ &= f(\mathbf{x}) + \epsilon \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \partial_{\mathbf{z}} - \int_0^\epsilon \frac{1}{\sigma_\theta^2} \left(\|\partial_{\mathbf{z}}\|_2^2 - \frac{1}{D_{\text{aux}}} \sum_{i=1}^{D_{\text{aux}}} (\langle \mathbf{1}, \partial_{\mathbf{z}} \rangle)^2 - (\langle \mathbf{z}_\theta, \partial_{\mathbf{z}} \rangle)^2 \mathbf{z}_\theta \right) \theta d\theta, \end{aligned}$$

where \mathbf{x}_θ represents $\mathbf{x} + \theta \partial_{\mathbf{z}}$, $\mathbf{z}_\theta = f(\mathbf{x}_\theta)$. The second step follows similar steps for computing $\frac{\partial \mathbf{z}}{\partial \mathbf{x}}$ in Equation (19). We avoid this computation since we only need to make sure that the second-order term is bounded. Furthermore, if $\epsilon \leq \mathcal{O}\left(\frac{\sigma}{\sqrt{D_{\text{aux}}}\|\partial_{\mathbf{z}}\|_2}\right)$, we can show the ℓ_2 -norm of the second-order term can be bounded by $\mathcal{O}(\epsilon^2 D_{\text{aux}}^{3/2} \sigma^{-2} \|\partial_{\mathbf{z}}\|_2^2)$. We avoid this computation as well.

Thus, from the above formulation, we have

$$\lim_{\epsilon \rightarrow 0} \frac{f(\mathbf{x} + \epsilon \partial_{\mathbf{z}}) - f(\mathbf{x})}{\epsilon} = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \partial_{\mathbf{z}} = \left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right)^\top \partial_{\mathbf{z}} = \partial_{\mathbf{x}}.$$

The pre-final step follows from Equation (19), where $\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \frac{1}{\sigma} \left(\mathbf{I} - \frac{1}{D_{\text{aux}}} \mathbf{1}\mathbf{1}^\top - \mathbf{z}\mathbf{z}^\top \right)$ can be shown to be symmetric. The final step follows from the gradient formulation in Equation (18). Including the error term, we have the final bound as

$$\left\| \frac{f(\mathbf{x} + \epsilon \partial_{\mathbf{z}}) - f(\mathbf{x})}{\epsilon} - \partial_{\mathbf{x}} \right\|_2 \leq \mathcal{O}(\epsilon D_{\text{aux}}^{3/2} \sigma^{-2} \|\partial_{\mathbf{z}}\|_2^2).$$

Using $\partial_{\mathbf{z}} = \gamma \odot \partial_{\mathbf{y}}$ and a Cauchy-Schwartz inequality gives the final bound. \square

G Activation layer

Definition G.1 (Auxiliary activation). For a continuous function $\sigma_{\text{act}} : \mathbb{R} \rightarrow \mathbb{R}$, an activation layer takes $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ as input and outputs $\mathbf{y} = \sigma_{\text{act}}(\mathbf{x})$ with $y_i = \sigma_{\text{act}}(x_i)$ for all $i \leq D_{\text{aux}}$.

In the discussions below, we consider an activation layer in the auxiliary model with activation function σ_{act} that takes in input sequence $\mathbf{x}_1, \dots, \mathbf{x}_{T_{\text{aux}}}$ and outputs $\mathbf{y}_1, \dots, \mathbf{y}_{T_{\text{aux}}}$, with $\mathbf{y}_t = \sigma_{\text{act}}(\mathbf{x}_t)$ for each $t \leq T_{\text{aux}}$. Since this involves a token-wise operation, we will present our constructed modules with a general token position t . Since no parameters of the auxiliary model are involved in this operation, the prefix tokens $\{\mathbf{v}_j\}$ contain 0 in the following modules.

TINT Activation Forward module The embedding \mathbf{e}_t contains \mathbf{x}_t in its first D_{aux} indices. We simply pass the embeddings into activation σ_{act} , which returns $\sigma_{\text{act}}(\mathbf{x}_t)$ in its first D_{aux} indices.

Auxiliary's backpropagation through activation With the definition in Definition G.1, the auxiliary's backpropagation takes in the loss gradient w.r.t. output ($\partial_{\mathbf{y}}$) and computes the loss gradient w.r.t. input ($\partial_{\mathbf{x}}$). We further assume that the derivative of σ_{act} is well-defined everywhere. This assumption includes non-differentiable activation functions with well-defined derivatives like *ReLU*.

Definition G.2 (Auxiliary activation backpropagation). For a continuous function $\sigma_{\text{act}} : \mathbb{R} \rightarrow \mathbb{R}$, with a well-defined derivative $\sigma'_{\text{act}}(x) = \partial \sigma_{\text{act}}(x) / \partial x$ for each $x \in \mathbb{R}$, the backpropagation takes $\partial_{\mathbf{y}}, \mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ as input and outputs

$$\partial_{\mathbf{x}} = \sigma'_{\text{act}}(\mathbf{x}) \odot \partial_{\mathbf{y}},$$

where $\sigma'_{\text{act}}(\mathbf{x}) \in \mathbb{R}^{D_{\text{aux}}}$ with $\sigma'_{\text{act}}(\mathbf{x})_i = \sigma'_{\text{act}}(x_i)$ at each $i \leq D_{\text{aux}}$.

Complexity of true backpropagation The above operation is computation heavy since it involves $\sigma'_{\text{act}}(\mathbf{x}) \odot \partial_{\mathbf{y}}$. As mentioned for the layer normalization module, the element-wise multiplication between $\sigma'_{\text{act}}(\mathbf{x})$ and $\partial_{\mathbf{y}}$ will require an MLP module following Lemma C.4. Furthermore, it involves changing the activation function in TINT in specific modules to σ'_{act} . To circumvent this, we instead turn to a first-order Taylor approximation.

Definition G.3 (Approximate Activation backpropagation). For a continuous function $\sigma_{\text{act}} : \mathbb{R} \rightarrow \mathbb{R}$ and a hyperparameter ϵ , the layer takes $\partial_{\mathbf{y}}, \mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ as input and outputs

$$\widehat{\partial_{\mathbf{x}}} = \frac{1}{\epsilon} (\sigma_{\text{act}}(\mathbf{x} + \epsilon \partial_{\mathbf{y}}) - \sigma_{\text{act}}(\mathbf{x})).$$

The following theorems show that under mild assumptions on the activation function and the input, gradient pair, the first-order gradient is a good approximation to the true gradient.

Theorem G.4. For any $\epsilon > 0$, $B_y, B_{\text{act}} > 0$, consider a second-order differentiable activation function $\sigma_{\text{act}} : \mathbb{R} \rightarrow \mathbb{R}$, with $\partial^2 \sigma_{\text{act}}(x) / \partial(x^2)$ bounded by B_{act} for each $x \in \mathbb{R}$. Then, for any input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ and gradient $\partial_{\mathbf{y}} \in \mathbb{R}^{D_{\text{aux}}}$ with $\|\partial_{\mathbf{y}}\|_2 \leq B_y$, the following holds true:

$$\|\partial_{\mathbf{x}} - \widehat{\partial_{\mathbf{x}}}\|_2 \leq \mathcal{O}(B_{\text{act}} B_y^2 \epsilon),$$

where $\partial_{\mathbf{x}}, \widehat{\partial_{\mathbf{x}}}$ have been defined using $\mathbf{x}, \partial_{\mathbf{y}}$, and ϵ in Definitions G.2 and G.3.

For ReLU activation, which is not second-order differentiable at 0, we instead bound the difference between $\partial_{\mathbf{x}}, \widehat{\partial_{\mathbf{x}}}$ by defining some form of alignment between input and gradient pair $\mathbf{x}, \partial_{\mathbf{y}}$.

Definition G.5 ((ϵ, ρ) -alignment). Input and gradient $\mathbf{x}, \partial_{\mathbf{y}} \in \mathbb{R}^{D_{\text{aux}}}$ are said to be (ϵ, ρ) -aligned, if there exist a set $C \subseteq [D_{\text{aux}}]$, with $|C| \geq (1 - \rho)D_{\text{aux}}$, such that for each i in C , $|x_i| > \epsilon |(\partial_{\mathbf{y}})_i|$.

ϵ controls the fraction of coordinates where $|x_i| \leq \epsilon |(\partial_{\mathbf{y}})_i|$. As $\epsilon \rightarrow 0$, $\rho \rightarrow 0$ as well for bounded gradients.

Example G.6. For any $B_{\min}, B_{\max} > 0$, all inputs \mathbf{x} that satisfy $\min_i |x_i| > B_{\min}$, and gradients $\partial_{\mathbf{y}}$ that satisfy $\max_j |(\partial_{\mathbf{y}})_j| \leq B_{\max}$, are $(B_{\min}/B_{\max}, 0)$ -aligned.

Theorem G.7. For any $\epsilon, \rho > 0$ and $B_y > 0$, for any input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ and gradient $\partial_{\mathbf{y}} \in \mathbb{R}^{D_{\text{aux}}}$, with $\|\partial_{\mathbf{y}}\|_{\infty} \leq B_y$, that are (ϵ, ρ) -aligned by Definition G.5,

$$\|\partial_{\mathbf{x}} - \widehat{\partial_{\mathbf{x}}}\|_2 \leq \mathcal{O}(B_y \sqrt{\rho D_{\text{aux}}}).$$

where $\partial_{\mathbf{x}}, \widehat{\partial_{\mathbf{x}}}$ have been defined using $\mathbf{x}, \partial_{\mathbf{y}}, \epsilon$ and $\sigma_{\text{act}} = \text{ReLU}$ in Definitions G.2 and G.3.

TINT Activation backpropagation module The input embeddings contain $\partial_{\mathbf{y}_t}$ in the first D_{aux} embeddings. With the requirement of the activation layer input for gradient, we copy \mathbf{x}_t from the Forward module at each position t . We set ϵ as a hyper-parameter and return $\widehat{\partial_{\mathbf{x}_t}}$ as the output of this module.

$\widehat{\partial_{\mathbf{x}_t}}$ will be computed using a single-layer MLP with activation σ_{act} as follows. The first linear layer of the MLP will be used to compute $\mathbf{x}_t + \epsilon \partial_{\mathbf{y}_t}$ and \mathbf{x}_t . After the activation σ_{act} , the embedding \mathbf{e}_t contains $\sigma_{\text{act}}(\mathbf{x}_t + \epsilon \partial_{\mathbf{y}_t})$ and $\sigma_{\text{act}}(\mathbf{x}_t)$. The final linear layer of the MLP will be used to compute $\frac{1}{\epsilon} (\sigma_{\text{act}}(\mathbf{x}_t + \epsilon \partial_{\mathbf{y}_t}) - \sigma_{\text{act}}(\mathbf{x}_t))$.

G.1 Proofs of theorems

We restate the theorems, before presenting their proofs for easy referencing.

Theorem G.4. *For any $\epsilon > 0$, $B_y, B_{act} > 0$, consider a second-order differentiable activation function $\sigma_{act} : \mathbb{R} \rightarrow \mathbb{R}$, with $\partial^2 \sigma_{act}(x)/\partial(x^2)$ bounded by B_{act} for each $x \in \mathbb{R}$. Then, for any input $\mathbf{x} \in \mathbb{R}^{D_{aux}}$ and gradient $\partial_{\mathbf{y}} \in \mathbb{R}^{D_{aux}}$ with $\|\partial_{\mathbf{y}}\|_2 \leq B_y$, the following holds true:*

$$\left\| \partial_{\mathbf{x}} - \widehat{\partial_{\mathbf{x}}} \right\|_2 \leq \mathcal{O}(B_{act} B_y^2 \epsilon),$$

where $\partial_{\mathbf{x}}, \widehat{\partial_{\mathbf{x}}}$ have been defined using $\mathbf{x}, \partial_{\mathbf{y}}$, and ϵ in Definitions G.2 and G.3.

Proof. The proof follows along the lines of Theorem F.1. Recall that given an input \mathbf{x} , the activation layer outputs $\mathbf{y} = \sigma_{act}(\mathbf{x})$, where the function σ_{act} is applied coordinate-wise on \mathbf{x} . Given input \mathbf{x} and the output gradient $\partial_{\mathbf{y}}$, the gradient w.r.t. the input is given by $\partial_{\mathbf{x}} = \sigma'_{act}(\mathbf{x}) \odot \partial_{\mathbf{y}}$, where the σ'_{act} function is also applied coordinate wise to \mathbf{x} . We defined $\widehat{\partial_{\mathbf{x}}}$ as an ϵ -approximate gradient, given by $\frac{1}{\epsilon}(\sigma_{act}(\mathbf{x} + \epsilon \partial_{\mathbf{y}}) - \sigma_{act}(\mathbf{x}))$. Since both σ_{act} and σ'_{act} are applied coordinate-wise, we can look at the coordinate-wise difference between $\partial_{\mathbf{x}}$ and $\widehat{\partial_{\mathbf{x}}}$.

Consider an arbitrary coordinate $i \leq D_{aux}$. Under the assumption that σ_{act} is second-order differentiable, we have

$$\begin{aligned} (\widehat{\partial_{\mathbf{x}}})_i &= \frac{1}{\epsilon} (\sigma_{act}(x_i + \epsilon(\partial_{\mathbf{y}})_i) - \sigma_{act}(x_i)) \\ &= \sigma'_{act}(x_i)(\partial_{\mathbf{y}})_i + \frac{1}{\epsilon} \int_{\theta=0}^{\epsilon} \frac{\partial^2 \sigma_{act}(x_{\theta})}{\partial x_{\theta}^2} (\partial_{\mathbf{y}})_i^2 \theta d\theta \\ &= \sigma'_{act}(x_i)(\partial_{\mathbf{y}})_i + \mathcal{O}(\epsilon B_{act} (\partial_{\mathbf{y}})_i^2), \end{aligned}$$

where x_{θ} represents $x_i + \theta(\partial_{\mathbf{y}})_i$ in the second step. In the final step, we utilize the upper bound assumption on $\frac{\partial^2 \sigma_{act}(x)}{\partial x^2}$.

Thus, $(\partial_{\mathbf{x}})_i - (\widehat{\partial_{\mathbf{x}}})_i = \mathcal{O}(\epsilon B_{act} (\partial_{\mathbf{y}})_i^2)$, and so

$$\left\| \partial_{\mathbf{x}} - \widehat{\partial_{\mathbf{x}}} \right\|_2 = \mathcal{O}(\epsilon B_{act} \sum_{i=1}^{D_{aux}} (\partial_{\mathbf{y}})_i^2) = \mathcal{O}(\epsilon B_{act} \|\partial_{\mathbf{y}}\|_2^2) \leq \mathcal{O}(\epsilon B_{act} B_y^2).$$

□

Example G.6. *For any $B_{min}, B_{max} > 0$, all inputs \mathbf{x} that satisfy $\min_i |x_i| > B_{min}$, and gradients $\partial_{\mathbf{y}}$ that satisfy $\max_j |(\partial_{\mathbf{y}})_j| \leq B_{max}$, are $(B_{min}/B_{max}, 0)$ -aligned.*

Proof. Recall the definition of (ϵ, ρ) -alignment from Definition G.5. Input and gradient $\mathbf{x}, \partial_{\mathbf{y}} \in \mathbb{R}^{D_{aux}}$ are said to be (ϵ, ρ) -aligned, if there exist a set $C \subseteq [D_{aux}]$, with $|C| \geq (1 - \rho)D_{aux}$, such that for each i in C , $|x_i| > \epsilon |(\partial_{\mathbf{y}})_i|$.

Consider an arbitrary coordinate $i \leq D_{aux}$. We have $|x_i| > \epsilon |(\partial_{\mathbf{y}})_i|$ for any $\epsilon < |x_i| / |(\partial_{\mathbf{y}})_i|$. Under the assumption that $|x_i| > B_{min}$, and $|(\partial_{\mathbf{y}})_i| \leq B_{max}$, a bound of B_{min}/B_{max} suffices. □

Theorem G.7. *For any $\epsilon, \rho > 0$ and $B_y > 0$, for any input $\mathbf{x} \in \mathbb{R}^{D_{aux}}$ and gradient $\partial_{\mathbf{y}} \in \mathbb{R}^{D_{aux}}$, with $\|\partial_{\mathbf{y}}\|_{\infty} \leq B_y$, that are (ϵ, ρ) -aligned by Definition G.5,*

$$\left\| \partial_{\mathbf{x}} - \widehat{\partial_{\mathbf{x}}} \right\|_2 \leq \mathcal{O}(B_y \sqrt{\rho D_{aux}}).$$

where $\partial_{\mathbf{x}}, \widehat{\partial_{\mathbf{x}}}$ have been defined using $\mathbf{x}, \partial_{\mathbf{y}}$, ϵ and $\sigma_{act} = \text{ReLU}$ in Definitions G.2 and G.3.

Proof. Recall that given an input \mathbf{x} , the activation layer outputs $\mathbf{y} = \sigma_{act}(\mathbf{x})$, where the function σ_{act} is applied coordinate-wise on \mathbf{x} . Given input \mathbf{x} and the output gradient $\partial_{\mathbf{y}}$, the gradient w.r.t. the input is given by $\partial_{\mathbf{x}} = \sigma'_{act}(\mathbf{x}) \odot \partial_{\mathbf{y}}$, where the σ'_{act} function is also applied coordinate wise to \mathbf{x} . We defined $\widehat{\partial_{\mathbf{x}}}$ as an ϵ -approximate gradient, given by $\frac{1}{\epsilon}(\sigma_{act}(\mathbf{x} + \epsilon \partial_{\mathbf{y}}) - \sigma_{act}(\mathbf{x}))$. Since both σ_{act} and

σ'_{act} are applied coordinate-wise, we can look at the coordinate-wise difference between $\partial_{\mathbf{x}}$ and $\widehat{\partial_{\mathbf{x}}}$. For ReLU activation, $\sigma'_{\text{act}}(x) = \text{sign}(x)$ for all $x \in \mathbb{R} \setminus \{0\}$, with $\sigma'_{\text{act}}(0) = 1$ to avoid ambiguity.

Going by the definition of (ϵ, ρ) -alignment of the input and gradient from Definition G.5, we have a set C with $|C| \geq (1 - \rho)D_{\text{aux}}$ such that for each $i \in D_{\text{aux}}$, $|x_i| > \epsilon |(\partial_{\mathbf{y}})_i|$. For all coordinates $i \in C$, we can then observe that $\text{sign}(x_i + \epsilon(\partial_{\mathbf{y}})_i) = \text{sign}(x_i)$, implying

$$\sigma_{\text{act}}(x_i + \epsilon(\partial_{\mathbf{y}})_i) - \sigma_{\text{act}}(x_i) = \epsilon(\partial_{\mathbf{y}})_i \sigma'_{\text{act}}(x_i) = \epsilon(\partial_{\mathbf{x}})_i$$

For coordinates $i \notin C$, we have three possible cases:

- $\text{sign}(x_i) = \text{sign}(x_i + \epsilon(\partial_{\mathbf{y}})_i)$: In this case, we can again show $\sigma_{\text{act}}(x_i + \epsilon(\partial_{\mathbf{y}})_i) - \sigma_{\text{act}}(x_i) = \epsilon(\partial_{\mathbf{y}})_i \sigma'_{\text{act}}(x_i) = \epsilon(\partial_{\mathbf{x}})_i$.
- $\text{sign}(x_i) = 0$, $\text{sign}(x_i + \epsilon(\partial_{\mathbf{y}})_i) = 1$: In this case, we have $\sigma'_{\text{act}}(x_i) = 0$, and so $(\partial_{\mathbf{x}})_i = 0$. Additionally, $\text{sign}((\partial_{\mathbf{y}})_i) = 1$, and so

$$|\sigma_{\text{act}}(x_i + \epsilon(\partial_{\mathbf{y}})_i) - \sigma_{\text{act}}(x_i) - \epsilon(\partial_{\mathbf{x}})_i| = |x_i + \epsilon(\partial_{\mathbf{y}})_i| \leq \epsilon |(\partial_{\mathbf{y}})_i|,$$

where in the final step, we use the fact that $x_i < 0$ and $|x_i| < \epsilon |(\partial_{\mathbf{y}})_i|$.

- $\text{sign}(x_i) = 1$, $\text{sign}(x_i + \epsilon(\partial_{\mathbf{y}})_i) = 0$: In this case, we have $\sigma'_{\text{act}}(x_i) = 1$, and so $(\partial_{\mathbf{x}})_i = (\partial_{\mathbf{y}})_i$. Additionally, $\text{sign}((\partial_{\mathbf{y}})_i) = 0$, and so

$$|\sigma_{\text{act}}(x_i + \epsilon(\partial_{\mathbf{y}})_i) - \sigma_{\text{act}}(x_i) - \epsilon(\partial_{\mathbf{x}})_i| = |-x_i - \epsilon(\partial_{\mathbf{y}})_i| \leq |\epsilon(\partial_{\mathbf{y}})_i|,$$

where in the final step, we use the fact that $x_i \geq 0$ and $|x_i| < \epsilon |(\partial_{\mathbf{y}})_i|$.

Thus, from the above discussion, we have

$$\begin{aligned} \|\partial_{\mathbf{x}} - \widehat{\partial_{\mathbf{x}}}\|_2 &= \frac{1}{\epsilon} \left(\sum_{i=1}^{D_{\text{aux}}} (\sigma_{\text{act}}(x_i + \epsilon(\partial_{\mathbf{y}})_i) - \sigma_{\text{act}}(x_i) - \epsilon(\partial_{\mathbf{x}})_i)^2 \right)^{1/2} \\ &= \frac{1}{\epsilon} \left(\sum_{i \notin C} (\sigma_{\text{act}}(x_i + \epsilon(\partial_{\mathbf{y}})_i) - \sigma_{\text{act}}(x_i) - \epsilon(\partial_{\mathbf{x}})_i)^2 \right)^{1/2} \\ &\leq \left(\sum_{i \notin C} (\partial_{\mathbf{y}})_i^2 \right)^{1/2} \leq \sqrt{\rho D_{\text{aux}}} \sqrt{\max_{i \notin C} (\partial_{\mathbf{y}})_i^2} \leq \sqrt{\rho D_{\text{aux}}} B_{\mathbf{y}}. \end{aligned}$$

The final step includes a simple Cauchy Schwartz inequality and the desired bound comes from the assumed bound on $\|\partial_{\mathbf{y}}\|_2$. \square

H Language model head

Additionally, we provide a description of the gradient computation for the loss function that involves the language model head. This computation entails performing a softmax operation over the entire vocabulary. If \mathcal{V} denotes the vocabulary set of the auxiliary model, and $\mathbf{E} \in \mathbb{R}^{|\mathcal{V}| \times D_{\text{aux}}}$ denotes the embedding matrix of the auxiliary model, we directly utilize the embedding matrix for the autoregressive loss in the TINT. Additionally, we do not update the embedding matrix of the auxiliary model; instead, we solely backpropagate the gradients through the language model head. Recent work in [21] has shown that keeping the embedding matrix fixed while updating the model can stabilize SGD. We demonstrate that the backpropagated gradients can be expressed as the combination of the language model head and a self-attention layer.

Definition H.1 (KL-loss gradient through auxiliary’s language model head). Given an embedding matrix $\mathbf{E} \in \mathbb{R}^{|\mathcal{V}| \times D_{\text{aux}}}$, the language model head takes in input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ and a target distribution $\mathbf{q} \in \mathbb{R}^{|\mathcal{V}|}$ and returns gradient $\partial_{\mathbf{x}} \in \mathbb{R}^{D_{\text{aux}}}$, with $\partial_{\mathbf{x}} = \mathbf{E}^\top (\text{softmax}(\mathbf{E}\mathbf{x}) - \mathbf{q})$.

In the autoregressive loss on a sequence of tokens, the target output distribution at any position is the next occurring token. If $\{\mathbf{x}_t^{\text{un}}\}_{t=1}^{T_{\text{aux}}}$ denote the uncontextualized embeddings of a sequence of tokens after encoding them via the embedding matrix, and $\{\mathbf{x}_t\}_{t=1}^{T_{\text{aux}}}$ denote their contextualized embeddings after passing through the auxiliary model, then the gradient $\partial_{\mathbf{x}_t}$ at any position t can be simplified as $\mathbf{E}^\top \text{softmax}(\mathbf{E}\mathbf{x}_t) - \mathbf{x}_{t+1}^{\text{un}}$. We illustrate the involved TINT module w.r.t. an arbitrary position t .

TINT autoregressive loss gradient module The current embedding e_t contains the contextualized embedding x_t in its first D_{aux} coordinates. Furthermore, e_t includes the uncontextualized embedding x_t^{un} , copied from the input layer using residual connections. The prefix tokens v_j are assigned a value of 0 and do not participate in the subsequent computations.

The loss computation can be decomposed into two sub-operations: (a) computing $y_t := E^\top \text{softmax}(Ex_t)$, and (b) calculating $\partial_{x_t} = y_t - x_{t+1}^{\text{un}}$.

For the first sub-operation, we use a feed-forward layer with softmax activation, with hidden and output weights E and E^\top respectively, that takes in the first D_{aux} of e_t and returns y_t in the first D_{aux} coordinates. We retain x_t^{un} using a residual connection.

The final sub-operation can be interpreted as a TINT self-attention layer. With e_t containing both y_t and x_t^{un} , we use a linear self-attention layer (Definition C.1) with two attention heads. The first attention head assigns an attention score of 1 to pairs $\{(t, t+1)\}_{t \leq T_{\text{aux}}-1}$, while assigning an attention score of 0 to the remaining pairs. At any position t , $-x_t^{\text{un}}$ is considered the value vector. The second attention head assigns an attention score of 1 to pairs $\{(t, t)\}_{t \leq T_{\text{aux}}}$, while assigning an attention score of 0 to the remaining pairs. At any position t , y_t is considered the value vector. The outputs of both attention heads are subsequently combined using a linear layer.

Remark H.2. We conducted experiments using mean-squared loss and Quad loss [35], which do not necessitate softmax computations for gradient computation. As an example, in the case of mean-squared loss, if our objective is to minimize $\frac{1}{2} \sum_{t=1}^T \|x_t - x_{t+1}^{\text{un}}\|^2$, the gradient can be computed as $\partial_{x_t} = x_t - x_{t+1}^{\text{un}}$. Similarly, in the case of Quad loss, the gradient is $\partial_{x_t} = \frac{1}{|V|} \sum_i e_i - x_{t+1}^{\text{un}}$. However, in all of our language model experiments, both gradients resulted in minimal improvement in perplexity compared to the auxiliary model. Therefore, we continue utilizing the standard KL loss for optimization.

Remark H.3. For ease of implementation in the codebase, we utilize a dedicated loss module that takes in y_t, x_{t+1}^{un} as input and directly computes $\partial_{x_t} = y_t - x_{t+1}^{\text{un}}$.

I Parameter sharing

Feed-forward layer of auxiliary model: In a standard auxiliary transformer, like GPT-2, the feed-forward layer is a token-wise operation that takes in an input $x \in \mathbb{R}^{D_{\text{aux}}}$ and returns $y = A\sigma(Wx)$, with $A \in \mathbb{R}^{D_{\text{aux}} \times 4D_{\text{aux}}}$ and $W \in \mathbb{R}^{4D_{\text{aux}} \times D_{\text{aux}}}$. A naive construction of the TINT to simulate its forward operation will have 2 Linear Forward modules (Section 2.4), separated by an activation. However, this requires $4 \times$ more prefix embeddings to represent the parameters, compared to other linear operations in the auxiliary transformer that use $\mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$ weight parameters.

To avoid this, we can instead break down the computation into 4 sub-feed-forward layers, each with its own parameters $\{\{W^i, A^i\}\}_{1 \leq i \leq 4}$. Here $\{W^i\}_{1 \leq i \leq 4}$ represent 4-shards of the rows of W , and $\{A^i\}_{1 \leq i \leq 4}$ represent 4-shards of the columns of A .

The forward, backward, and descent operations on these 4 sub-feed-forward layers can be effectively parallelized. For example, the forward operation of each layer can be simulated by a single TINT module, consisting of two Linear Forward modules and activation, changing only the prefix embeddings to correspond to $\{\{W^i, A^i\}\}_{1 \leq i \leq 4}$.

J Additional modules

We describe the forward, backward, and decent update operations of additional modules, used in different model families, like LLaMA [41] and BLOOM [36]. We discuss the simulation of these modules, using similar TINT modules.

J.1 Root mean square normalization (RMSnorm)

The operation of RMSnorm [52] is very similar to layer normalization.

Definition J.1 (RMSnorm). For an arbitrary dimension d , define a normalization function $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that performs $f(x) = x / \text{RMS}(x)$, where $\text{RMS}(x) = (\sum_{i=1}^d x_i^2)^{1/2}$. Then, RMSnorm with

parameters $\gamma, \mathbf{b} \in \mathbb{R}^{D_{\text{aux}}}$ takes as input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ and outputs $\mathbf{y} \in \mathbb{R}^{D_{\text{aux}}}$, which is computed as $\mathbf{z} = f(\mathbf{x}), \mathbf{y} = \gamma \odot \mathbf{z} + \mathbf{b}$.

The extreme similarity between RMSnorm and layer normalization (Definition A.1) helps us create similar TINT modules as described in Appendix F, where instead of Group normalization layers, we use Group RMSnorm layers described below.

Definition J.2 (TINT D_{aux} -Group RMSnorm). For an arbitrary dimension d , define a normalization function $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that performs $f(\mathbf{x}) = \mathbf{x} / \text{RMS}(\mathbf{x})$, where $\text{RMS}(\mathbf{x}) = (\sum_{i=1}^d x_i^2)^{1/2}$. Then, D_{aux} -Group RMSnorm with parameters $\gamma^{\text{TINT}}, \mathbf{b}^{\text{TINT}} \in \mathbb{R}^{D_{\text{aux}}}$ takes as input $\mathbf{x} \in \mathbb{R}^{D_{\text{sim}}}$ and outputs $\mathbf{y} = \text{VECTORIZE}(\{\mathbf{y}^h \in \mathbb{R}^{D_{\text{aux}}}\}_{h \leq \lfloor D_{\text{sim}}/D_{\text{aux}} \rfloor})$, with

$$\mathbf{y}^h = \gamma^{\text{TINT}} \odot f(\mathbf{x}^h) + \mathbf{b}^{\text{TINT}},$$

where $\mathbf{x}^h = \text{SPLIT}_{\lfloor D_{\text{sim}}/D_{\text{aux}} \rfloor}(\mathbf{x})_h$.

J.2 Attention variants

In order to incorporate additional attention variants, e.g. Attention with Linear Biases (ALiBi) [31], and rotary position embeddings [40], we can change the definition of softmax attention layer in Definition C.1 likewise.

J.3 Gated linear units (GLUs)

We describe the operations of GLUs [38] using similar GLU units available to the TINT.

Definition J.3. For parameters $\mathbf{W}, \mathbf{V}, \mathbf{W}^o \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$, and biases $\mathbf{b}_W, \mathbf{b}_V, \mathbf{b}_{W^o} \in \mathbb{R}^{D_{\text{aux}}}$, a GLU layer with activation $\sigma_{\text{act}} : \mathbb{R} \rightarrow \mathbb{R}$, takes input $\mathbf{x} \in \mathbb{R}^{D_{\text{aux}}}$ and outputs $\hat{\mathbf{y}} \in \mathbb{R}^{D_{\text{aux}}}$, with

$$\mathbf{y} = (\mathbf{W}\mathbf{x} + \mathbf{b}_W) \odot \sigma_{\text{act}}(\mathbf{V}\mathbf{x} + \mathbf{b}_V); \quad \hat{\mathbf{y}} = \mathbf{W}^o \mathbf{y} + \mathbf{b}_{W^o}.$$

Typical GLUs have $8/3 \times D_{\text{aux}}$ as a hidden dimension (i.e. the dimension of \mathbf{y}). We can use similar parameter-sharing techniques discussed for feed-forward layers (Appendix I) with the TINT modules presented here. Furthermore, since $\hat{\mathbf{y}}$ can be expressed as a combination of the gated operation and a linear operation, we focus on the computation of \mathbf{y} here.

For the discussion below, we consider a GLU (without the output linear layer) in the auxiliary model, with parameters $\mathbf{W}, \mathbf{V}, \mathbf{b}_W, \mathbf{b}_V$, that takes in input sequence $\mathbf{x}_1, \dots, \mathbf{x}_T$ and outputs $\mathbf{y}_1, \dots, \mathbf{y}_T$, with $\mathbf{y}_t = (\mathbf{W}\mathbf{x}_t + \mathbf{b}_W) \odot \sigma_{\text{act}}(\mathbf{V}\mathbf{x}_t + \mathbf{b}_V)$ for each $t \leq T_{\text{sim}}$. Since this involves a token-wise operation, we will present our constructed modules with a general token position t and the prefix tokens $\{\mathbf{v}_j\}$.

TINT GLU Forward module The embedding \mathbf{e}_t contains \mathbf{x}_t in its first D_{aux} coordinates. The output \mathbf{y}_t can be computed using three sub-operations: (a) linear operation for $\mathbf{W}\mathbf{x}_t + \mathbf{b}_W$, (b) linear operation for $\mathbf{V}\mathbf{x}_t + \mathbf{b}_V$, and (c) gate operation to get $(\mathbf{W}\mathbf{x}_t + \mathbf{b}_W) \odot \sigma_{\text{act}}(\mathbf{V}\mathbf{x}_t + \mathbf{b}_V)$.

We use three TINT modules, representing each sub-operation.

- (a) $\mathbf{W}\mathbf{x}_t + \mathbf{b}_W$ is a linear operation, hence we can use a TINT Linear Forward module (Appendix D) with the current embedding \mathbf{e}_t and $\{\mathbf{v}_j\}$ containing \mathbf{W}, \mathbf{b}_W to get embedding $\tilde{\mathbf{e}}_t$ containing $\mathbf{W}\mathbf{x}_t + \mathbf{b}_W$ in its first D_{aux} coordinates.
- (b) $\mathbf{V}\mathbf{x}_t + \mathbf{b}_V$ is a linear operation, hence we can similarly use a TINT Linear Forward module (Appendix D) with the embedding \mathbf{e}_t and $\{\mathbf{v}_j\}$ containing $\mathbf{W}_V, \mathbf{b}_V$ to get embedding $\hat{\mathbf{e}}_t$ containing $\mathbf{V}\mathbf{x}_t + \mathbf{b}_V$ in its first D_{aux} coordinates.
 $\hat{\mathbf{e}}_t$ and $\tilde{\mathbf{e}}_t$ are now combined to get an embedding \mathbf{e}_t that contains $\mathbf{W}\mathbf{x}_t + \mathbf{b}_W, \mathbf{V}\mathbf{x}_t + \mathbf{b}_V$ in its first $2D_{\text{aux}}$ coordinates.
- (c) Finally, we can use a TINT GLU layer that can carry out the elementwise multiplication of $\mathbf{W}\mathbf{x}_t + \mathbf{b}_W, \sigma_{\text{act}}(\mathbf{V}\mathbf{x}_t + \mathbf{b}_V)$ to get \mathbf{y}_t in the first D_{aux} coordinates.

Parameter Sharing: Since (a) and (b) involve a Linear Forward module, we can additionally leverage parameter sharing to apply a single Linear Forward module for each of the two computations, changing only the prefix embeddings to correspond to \mathbf{W}, \mathbf{b}_W , or $\mathbf{W}_V, \mathbf{b}_V$.

Auxiliary GLU backpropagation For the GLU layer defined in Definition J.3, the backpropagation layer takes in the loss gradient w.r.t. output ($\partial_{\mathbf{y}}$) and computes the loss gradient w.r.t. input ($\partial_{\mathbf{x}}$).

Definition J.4 (Auxiliary GLU backpropagation). For the weights $\mathbf{W}, \mathbf{V} \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$, the backpropagation layer takes $\partial_{\mathbf{y}} \in \mathbb{R}^{D_{\text{aux}}}$ as input and outputs $\partial_{\mathbf{x}} \in \mathbb{R}^{D_{\text{aux}}}$, with $\partial_{\mathbf{x}} = \mathbf{W}^\top \widehat{\partial_{\mathbf{x}}} + \mathbf{V}^\top \widetilde{\partial_{\mathbf{x}}}$, where

$$\widehat{\partial_{\mathbf{x}}} = \partial_{\mathbf{y}} \odot \sigma_{\text{act}}(\mathbf{V}\mathbf{x} + \mathbf{b}_V); \quad \widetilde{\partial_{\mathbf{x}}} = \sigma'_{\text{act}}(\mathbf{V}\mathbf{x} + \mathbf{b}_V) \odot \partial_{\mathbf{y}} \odot (\mathbf{W}\mathbf{x} + \mathbf{b}_W).$$

A direct computation of $\widetilde{\partial_{\mathbf{x}}}$ involves changing the activation function to σ'_{act} . Following a similar strategy for backpropagation through an activation layer (Appendix G), we instead use a first-order Taylor expansion to approximate $\widetilde{\partial_{\mathbf{x}}}$.

Definition J.5 (Auxiliary GLU approximate backpropagation). For a hyper-parameter $\epsilon > 0$, for the weights $\mathbf{W}, \mathbf{V} \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$, the approximate backpropagation layer takes $\partial_{\mathbf{y}} \in \mathbb{R}^{D_{\text{aux}}}$ as input and outputs $\overline{\partial_{\mathbf{x}}} \in \mathbb{R}^{D_{\text{aux}}}$, with $\overline{\partial_{\mathbf{x}}} = \mathbf{W}^\top \widehat{\partial_{\mathbf{x}}} + \mathbf{V}^\top \widehat{\widetilde{\partial_{\mathbf{x}}}}$, where

$$\begin{aligned} \widehat{\partial_{\mathbf{x}}} &= \partial_{\mathbf{y}} \odot \sigma_{\text{act}}(\mathbf{V}\mathbf{x} + \mathbf{b}_V) \\ \widehat{\widetilde{\partial_{\mathbf{x}}}} &= \sigma_{\text{act}}(\mathbf{V}\mathbf{x} + \mathbf{b}_V + \epsilon \partial_{\mathbf{y}}) \odot \frac{1}{\epsilon}(\mathbf{W}\mathbf{x} + \mathbf{b}_W) - \sigma_{\text{act}}(\mathbf{V}\mathbf{x} + \mathbf{b}_V) \odot \frac{1}{\epsilon}(\mathbf{W}\mathbf{x} + \mathbf{b}_W). \end{aligned}$$

TINT GLU backpropagation module The current embedding contains $\partial_{\mathbf{y}_t}$ in its first D_{aux} coordinates. Furthermore, since we need $\mathbf{W}\mathbf{x}_t + \mathbf{b}_W$ and $\mathbf{V}\mathbf{x}_t + \mathbf{b}_V$ in the gradient computations, we copy them from the Forward module using residual connections. We discuss the computation of $\mathbf{W}^\top \widehat{\partial_{\mathbf{x}_t}}$ and $\mathbf{V}^\top \widehat{\widetilde{\partial_{\mathbf{x}_t}}}$ as separate sub-modules acting on the same embedding \mathbf{e}_t in parallel.

1. The computation of $\mathbf{W}^\top \widehat{\partial_{\mathbf{x}_t}}$ involves two sub-operations: (a) gate operation to get $\widehat{\mathbf{x}_t} := \partial_{\mathbf{y}_t} \odot \sigma_{\text{act}}(\mathbf{V}\mathbf{x}_t + \mathbf{b}_V)$, and (b) linear backward operation to get $\mathbf{W}^\top \widehat{\mathbf{x}_t}$. Since for this operation, we require \mathbf{W} , we copy the contents of the prefix embeddings containing \mathbf{W}, \mathbf{b}_W from the Forward module.
 - (a) Since the current embedding \mathbf{e}_t contains both $\partial_{\mathbf{y}_t}$ and $\mathbf{W}\mathbf{x}_t + \mathbf{b}_W$, we can use a TINT GLU layer to get an embedding $\widehat{\mathbf{e}}_t^{(1)}$ that contains $\widehat{\partial_{\mathbf{x}_t}}$.
 - (b) The final linear backward operation can be performed by using a TINT Linear backpropagation module (Appendix D) with the embeddings $\widehat{\mathbf{e}}_t^{(1)}$ and the prefix embeddings. The final embedding $\widehat{\mathbf{e}}_t$ contains $\mathbf{W}^\top \widehat{\mathbf{x}_t}$ in the first D_{aux} coordinates.
2. The computation of $\mathbf{V}^\top \widehat{\widetilde{\partial_{\mathbf{x}_t}}}$ involves four sub-operations: (a) gate operation to get $\frac{1}{\epsilon}(\mathbf{W}\mathbf{x}_t + \mathbf{b}_W) \odot \sigma_{\text{act}}(\mathbf{V}\mathbf{x}_t + \mathbf{b}_V + \epsilon \partial_{\mathbf{y}_t})$, (b) gate operation to get $\frac{1}{\epsilon}(\mathbf{W}\mathbf{x}_t + \mathbf{b}_W) \odot \sigma_{\text{act}}(\mathbf{V}\mathbf{x}_t + \mathbf{b}_V)$, (c) a linear layer to compute $\widehat{\mathbf{x}_t}$, (c) linear backward operation to get $\mathbf{V}^\top \widehat{\mathbf{x}_t}$. Since for this operation, we require \mathbf{V} , we copy the contents of the prefix embeddings containing \mathbf{V}, \mathbf{b}_V from the Forward module.
 - (a) Since the current embedding \mathbf{e}_t contains $\partial_{\mathbf{y}_t}, \mathbf{V}\mathbf{x}_t + \mathbf{b}_V$ and $\mathbf{W}\mathbf{x}_t + \mathbf{b}_W$, we can use two TINT GLU layers to get an embedding $\widetilde{\mathbf{e}}_t^{(1)}$ that contains both $\frac{1}{\epsilon}(\mathbf{W}\mathbf{x}_t + \mathbf{b}_W) \odot \sigma_{\text{act}}(\mathbf{V}\mathbf{x}_t + \mathbf{b}_V + \epsilon \partial_{\mathbf{y}_t})$ and $\frac{1}{\epsilon}(\mathbf{W}\mathbf{x}_t + \mathbf{b}_W) \odot \sigma_{\text{act}}(\mathbf{V}\mathbf{x}_t + \mathbf{b}_V)$.
 - (b) A linear later on $\widetilde{\mathbf{e}}_t^{(1)}$ can then return an embedding $\widetilde{\mathbf{e}}_t^{(2)}$ containing $\widehat{\mathbf{x}_t}$ in the first D_{aux} coordinates.
 - (c) The final operation can be performed by using a TINT Linear backpropagation module (Appendix D) with the embeddings $\widetilde{\mathbf{e}}_t^{(2)}$ and the prefix embeddings containing \mathbf{V}, \mathbf{b}_V . The final embedding $\widetilde{\mathbf{e}}_t$ contains $\mathbf{V}^\top \widehat{\mathbf{x}_t}$ in the first D_{aux} coordinates.

After the two parallel computations, we can sum up $\widehat{\mathbf{e}}_t$ and $\widetilde{\mathbf{e}}_t$ to get an embedding \mathbf{e}_t containing $\overline{\partial_{\mathbf{x}_t}}$ (Definition J.5) in the first D_{aux} coordinates.

Auxiliary GLU descent Finally, the auxiliary’s descent updates the weight and the bias parameters using a batch of inputs $\{\mathbf{x}_t\}_{t \leq T}$ and the loss gradient w.r.t. the corresponding outputs $\{\partial_{\mathbf{y}_t}\}_{t \leq T}$.

Definition J.6 (Auxiliary GLU descent). For weights $\mathbf{W}, \mathbf{V} \in \mathbb{R}^{D_{\text{aux}} \times D_{\text{aux}}}$ and bias $\mathbf{b}_W, \mathbf{b}_V \in \mathbb{R}^{D_{\text{aux}}}$, the linear descent layer takes in a batch of inputs $\{\mathbf{x}_t \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and gradients $\{\partial_{\mathbf{y}_t} \in \mathbb{R}^{D_{\text{aux}}}\}_{t \leq T_{\text{aux}}}$ and updates the parameters as follows:

$$\begin{aligned} \mathbf{W} &\leftarrow \mathbf{W} - \eta \sum_{t \leq T_{\text{aux}}} \widehat{\partial_{\mathbf{x}_t}} \mathbf{x}_t^\top; & \mathbf{b}_W &\leftarrow \mathbf{b}_W - \eta \sum_{t \leq T_{\text{aux}}} \widehat{\partial_{\mathbf{x}_t}}, \\ \mathbf{V} &\leftarrow \mathbf{V} - \eta \sum_{t \leq T_{\text{aux}}} \widetilde{\partial_{\mathbf{x}_t}} \mathbf{x}_t^\top; & \mathbf{b}_V &\leftarrow \mathbf{b}_V - \eta \sum_{t \leq T_{\text{aux}}} \widetilde{\partial_{\mathbf{x}_t}}, \end{aligned}$$

where $\widehat{\partial_{\mathbf{x}_t}}$ and $\widetilde{\partial_{\mathbf{x}_t}}$ have been computed as Definition J.4.

Due to similar concerns as gradient backpropagation, we instead use $\widehat{\partial_{\mathbf{x}_t}}$ (Definition J.5) in place of $\widetilde{\partial_{\mathbf{x}_t}}$ for each $t \leq T_{\text{aux}}$ to update \mathbf{V}, \mathbf{b}_V .

TINT GLU descent module We discuss the two descent operations separately.

1. Update of \mathbf{W}, \mathbf{b}_W : We start with the embeddings $\widehat{\mathbf{e}}_t^{(1)}$ from the backpropagation module, that contain $\widehat{\partial_{\mathbf{x}_t}}$ in the first D_{aux} coordinates.
For the update, we additionally require the input to the auxiliary GLU layer under consideration, and hence we copy \mathbf{x}_t from the Forward module using residual connections. Furthermore, we copy the contents of the prefix embeddings that contain \mathbf{W}, \mathbf{b}_W from the Forward module.
With both $\widehat{\partial_{\mathbf{x}_t}}$ and \mathbf{x}_t in the embeddings, the necessary operation turns out to be the descent update of a linear layer with parameters \mathbf{W}, \mathbf{b}_W . That implies, we can call a TINT Linear descent module (Appendix D) on the current embeddings and prefix embeddings to get the desired update.
2. We start with the embeddings $\widetilde{\mathbf{e}}_t^{(2)}$ from the backpropagation module, that contain $\widetilde{\partial_{\mathbf{x}_t}}$ in the first D_{aux} coordinates.
For the update, we additionally require the input to the auxiliary GLU layer under consideration, and hence we copy \mathbf{x}_t from the forward module using residual connections. Furthermore, we copy the contents of the prefix embeddings that contain \mathbf{V}, \mathbf{b}_V from the Forward module.
With both $\widetilde{\partial_{\mathbf{x}_t}}$ and \mathbf{x}_t in the embeddings, the necessary operation turns out to be the descent update of a linear layer with parameters \mathbf{V}, \mathbf{b}_V . That implies we can call a TINT Linear descent module on the current embeddings and prefix embeddings to get the desired update.

Parameter sharing: Since both the descent updates involve a Linear descent module, we can additionally leverage parameter sharing to apply a single TINT Linear descent module for each of the two computations, changing the input to correspond to $\{\widehat{\mathbf{e}}_t^{(1)}\}$ and prefix to correspond to \mathbf{W}, \mathbf{b}_W , or the input to correspond to $\{\widetilde{\mathbf{e}}_t^{(2)}\}$ and prefix to correspond to \mathbf{V}, \mathbf{b}_V respectively.

K Construction of other variants of pre-trained models

Table 3: Number of parameters of TINT for the forward, backward, and gradient update operations on various modules. For simplicity, we have ignored biases in the following computation. We set $H_{\text{sim}} = 12$ for OPT-125M and $H_{\text{sim}} = 16$ for the other models, $D_{\text{sim}} = 4D_{\text{aux}}$ for all the models, and $T_{\text{sim}} = T_{\text{aux}} + K$, with $T_{\text{aux}} = 2048$ for OPT models, and $K = D_{\text{aux}}/4$. $Q = 4Q_{\text{split}} + 3T_{\text{sim}}D_{\text{sim}}/H_{\text{sim}}$, where $Q_{\text{split}} = \frac{1}{H_{\text{sim}}}(D_{\text{sim}})^2 + H_{\text{sim}}D_{\text{sim}}$, denotes the number of parameters in a TINT Linear Forward module (Section 2.4).

Module Name	Module Size			
	Forward	Backward	Descent	Total
Linear layer	Q	Q	Q	$3Q$
Layer norms	Q	$Q + 2D_{\text{sim}}H_{\text{sim}}$	Q	$3Q + 2D_{\text{sim}}H_{\text{sim}}$
Self-Attention	$2Q$	$2Q$	$2Q$	$6Q$
Activation	Q_{split}	$2D_{\text{sim}}H_{\text{sim}}$	0	$Q_{\text{split}} + 2D_{\text{sim}}H_{\text{sim}}$
Self-Attention block	$4Q$	$4Q + 2D_{\text{sim}}H_{\text{sim}}$	$4Q$	$12Q + 2D_{\text{sim}}H_{\text{sim}}$
Feed-forward block	$3Q$	$3Q + 4D_{\text{sim}}H_{\text{sim}}$	$3Q$	$9Q + 4D_{\text{sim}}H_{\text{sim}}$
Transformer block	$7Q$	$7Q + 6D_{\text{sim}}H_{\text{sim}}$	$7Q$	$21Q + 6D_{\text{sim}}H_{\text{sim}}$
Transformer	$7QL + LQ_{\text{split}}$	$(7Q + 6D_{\text{sim}}H_{\text{sim}})L$	$7QL$	$(21Q + 6D_{\text{sim}}H_{\text{sim}})L$
OPT-125M	0.4B	0.4B	0.4B	1.2B
OPT-350M	1.2B	1.1B	1.1B	3.4B
OPT-1.3B	3.7B	3.6B	3.5B	10.8B
OPT-2.7B	7.4B	7.2B	7.2B	21.8B

Though we only conduct experiments on an OPT-125M model, our construction is generally applicable to diverse variants of pre-trained language models. Table 3 highlights many types of modules and the required size and computation for each. The size of a constructed model is influenced by various factors, including the number of layers, and embedding dimension in the auxiliary.

L Experiments

L.1 Experimental Setup

We introduce dynamic evaluation [20], which updates the model with a segment of the input and evaluates the updated model on the rest of the input.

Definition L.1 (Dynamic Evaluation [20]). For a sequence $s \in \mathcal{S}$ of T tokens s_1, \dots, s_T and hyperparameters i and η , *dynamic evaluation* of any model f with parameters $\theta \in \Theta$ is defined as follows. Let $\mathcal{L} : \{\mathcal{S}, \Theta\} \rightarrow \mathbb{R}$ be the cross-entropy language modeling loss. Compute

$$\theta' = \theta - \eta \nabla_{\theta} \mathcal{L}((s_1, \dots, s_i), \theta).$$

Then, the dynamic evaluation loss is $\mathcal{L}((s_{i+1}, \dots, s_T), \theta')$. For in-context learning, the first i examples are used to update the parameters.

For example, suppose we perform dynamic evaluation with $i = 5$ using the cross-entropy objective on **Machine learning is a useful tool for solving problems**. The **red** part is used to update the pre-trained model, and the **brown** part is used for evaluation. This operation is equivalent to TINT inference, which implicitly updates a model on the **red** portion during the forward pass and then evaluates the updated model on the **brown** portion. Therefore, we use dynamic evaluation as a baseline for TINT: matching dynamic evaluation performance provides strong evidence that TINT closely approximates explicitly fine-tuning the auxiliary model.

Tasks: We perform language modeling experiments on Wikitext-103 [25] and evaluate on 7 downstream tasks in zero-shot and few-shot settings: SST-2 [39], MR [28], CR [18], MPQA [48], Amazon Polarity [54], AGNews [54], and Subj [29].

Model: We compare a TINT model that tunes an OPT-125M pre-trained model internally to dynamic evaluation of OPT-125M and standard evaluation of OPT-1.3B.⁵ Given a downstream task input

⁵Our construction is generally applicable to diverse variants of pre-trained language models (Appendix K).

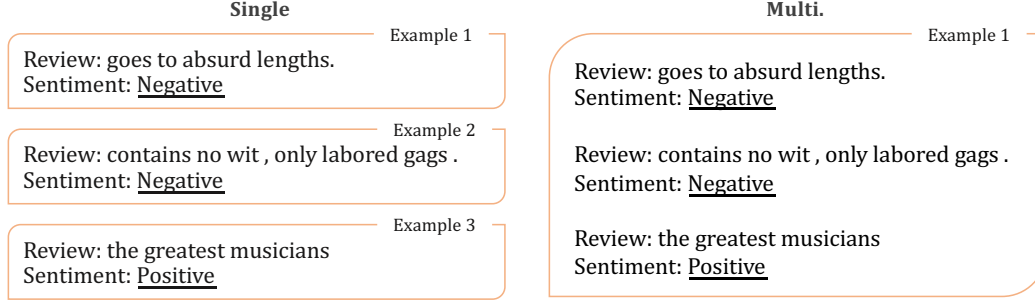


Figure 9: This illustration showcases different settings in few-shot learning ($k = 3$) using TINT. The **Single** mode (left) has one example for each input, and the auxiliary model is updated with a batch of inputs. The **Multi.** mode (right) concatenates all examples to form a single input. For **Label loss**, only underlined label words are used for internal training, while **full context loss** includes all tokens.

(e.g., a movie review), the model’s predicted label is computed as follows. First, we design a simple task-specific prompt (e.g., “Sentiment:”) and select label words c_1, \dots, c_n to serve as surrogates for each class (e.g., “positive” and “negative”). Then, we provide the input along with the prompt to the model, and the label word assigned the highest probability is treated as the model’s prediction. If using calibration, then the probabilities are normalized using just the prompt as input.⁶ Mathematically, the model’s prediction is computed as

$$\text{No calibration: } \arg \max_{c_i} \Pr[c_i \mid \text{input, prompt}] \quad \text{Calibration: } \arg \max_{c_i} \frac{\Pr[c_i \mid \text{input, prompt}]}{\Pr[c_i \mid \text{prompt}]}$$

This is a widely used calibration technique [17] for prompting language models.

Settings (Figure 9): We explore the following settings in downstream tasks: 1) Single and Multi.: We finetune the auxiliary model using either separate single examples or concatenated examples within each input; 2) Label loss and full-context loss: We finetune on the loss either from only label words or the entire context (Figure 9). We evaluate both zero-shot and few-shot settings, using the context of the evaluation example and 32 training examples for internal learning respectively.

L.2 Additional Experiments

Computing environment: All the experiments are conducted on a single A100 80G GPU.

Hyperparameters: In the few-shot setting, we employ three different random seeds to select distinct sets of training examples. Grid search is performed for each seed to determine the optimal learning rate for both constructed models and dynamic evaluation. The learning rates considered for the learning rate hyperparameter in the descent update operations in TINT are $1e-3$, $1e-4$, $1e-5$.⁷ Additionally, we explore various layer-step combinations to allocate a fixed budget for one full forward pass. Specifically, we update the top 3 layers for 4 steps, the top 6 layers for 3 steps, or 12 layers for 1 step.

Results of different settings. Table 4 displays the results of few-shot learning with calibration across various settings, encompassing different loss types, input formats, and layer-step configurations. Our analysis reveals that employing a label-only loss, utilizing a single-example input format, and updating all layers of the internal model for a single step yield the most favorable average result. The performance of the multi-example format is disadvantaged when dealing with tasks of long sequences such as Amazon Polarity. In general, we observe that calibrated results tend to be more consistent and stable.

⁶Calibration is not applied to the language modeling evaluation.

⁷When utilizing the full-context loss, the learning rates considered are $1e-5$, $1e-6$, and $1e-7$ due to gradient summations in TINT.

Table 4: Few-shot ($k = 32$) results with different loss types, input formats, and layer-step configurations with a fixed compute budget, with calibration.

Loss Type	Format	Layer	Step	Subj	AGNews	SST2	CR	MR	MPQA	Amazon	Avg.
Label	Single	12	1	66.0 _(1.9)	64.7 _(0.2)	68.7 _(1.3)	69.0 _(0.7)	63.7 _(0.2)	82.8 _(0.5)	73.7 _(0.6)	69.8 _(0.1)
	Single	6	2	62.7 _(0.2)	66.3 _(0.2)	68.3 _(6.1)	67.2 _(0.2)	61.8 _(1.6)	81.0 _(3.6)	74.3 _(0.5)	68.8 _(1.4)
	Single	3	4	63.5 _(0.0)	67.2 _(0.8)	62.5 _(0.4)	68.7 _(1.4)	61.7 _(0.6)	76.8 _(3.3)	75.2 _(0.8)	67.9 _(0.8)
	Multi.	12	1	83.2 _(2.5)	43.7 _(6.6)	60.7 _(5.7)	70.3 _(6.1)	62.8 _(8.9)	84.2 _(1.6)	66.3 _(12.3)	67.3 _(0.9)
	Multi.	6	2	83.5 _(2.9)	43.2 _(8.4)	52.0 _(1.5)	70.5 _(6.0)	58.5 _(11.3)	82.0 _(0.4)	55.8 _(7.6)	63.6 _(2.7)
	Multi.	3	4	84.0 _(2.3)	42.3 _(8.4)	51.5 _(1.8)	68.2 _(4.6)	58.5 _(12.0)	80.2 _(2.1)	58.5 _(7.9)	63.3 _(3.0)
Full-context	Single	12	1	64.5 _(0.4)	65.8 _(0.2)	63.2 _(0.9)	67.3 _(0.5)	60.8 _(1.4)	73.5 _(0.8)	75.0 _(0.4)	67.2 _(0.1)
	Single	6	2	66.7 _(2.0)	66.0 _(0.4)	62.7 _(0.6)	70.5 _(2.1)	59.7 _(0.9)	77.7 _(2.2)	76.0 _(0.0)	68.5 _(0.4)
	Single	3	4	64.0 _(0.0)	65.8 _(0.6)	65.0 _(1.9)	67.3 _(0.2)	59.5 _(0.4)	74.2 _(1.3)	77.0 _(1.9)	67.5 _(0.8)
	Multi.	12	1	83.8 _(2.9)	41.0 _(10.6)	51.2 _(0.8)	68.0 _(4.5)	58.3 _(11.1)	79.0 _(3.6)	56.0 _(8.1)	62.5 _(2.8)
	Multi.	6	2	85.3 _(1.9)	41.2 _(10.7)	51.2 _(1.3)	67.7 _(4.5)	57.7 _(10.8)	79.2 _(3.7)	55.8 _(7.9)	62.6 _(2.6)
	Multi.	3	4	83.3 _(2.5)	41.7 _(11.3)	51.0 _(1.1)	68.2 _(4.7)	57.7 _(10.8)	79.0 _(3.2)	56.0 _(8.1)	62.4 _(2.8)

Table 5: Few-shot ($k = 32$) results with different loss types, input formats, and layer-step configurations with a fixed compute budget, without calibration.

Loss Type	Format	Layer	Step	Subj	AGNews	SST2	CR	MR	MPQA	Amazon	Avg.
Label	Single	12	1	63.3 _(0.2)	65.7 _(0.2)	71.3 _(0.6)	65.0 _(1.4)	70.7 _(0.9)	65.0 _(0.0)	76.7 _(0.2)	68.2 _(0.1)
	Single	6	2	63.5 _(0.0)	65.2 _(0.5)	73.3 _(1.3)	68.5 _(3.7)	71.3 _(0.2)	66.0 _(0.0)	77.5 _(0.4)	69.3 _(0.3)
	Single	3	4	64.2 _(0.2)	66.5 _(1.1)	73.2 _(0.6)	75.7 _(0.5)	72.0 _(0.0)	83.2 _(1.0)	78.0 _(0.4)	73.2 _(0.1)
	Multi.	12	1	64.5 _(7.8)	35.5 _(7.4)	56.8 _(9.7)	63.0 _(6.7)	58.7 _(8.9)	75.2 _(10.8)	62.2 _(8.3)	59.4 _(0.6)
	Multi.	6	2	77.7 _(7.0)	35.5 _(7.4)	57.0 _(9.9)	60.0 _(6.3)	52.3 _(2.1)	58.5 _(6.1)	55.8 _(7.9)	56.7 _(2.6)
	Multi.	3	4	67.5 _(11.5)	38.5 _(8.2)	55.3 _(5.2)	67.0 _(3.5)	61.0 _(8.0)	65.2 _(11.2)	62.5 _(8.9)	59.6 _(1.3)
Full-context	Single	12	1	65.5 _(1.1)	66.5 _(0.0)	70.7 _(0.2)	64.8 _(0.5)	72.0 _(1.4)	67.0 _(0.0)	76.5 _(0.0)	69.0 _(0.3)
	Single	6	2	64.7 _(0.6)	66.2 _(0.2)	71.2 _(0.2)	65.3 _(0.6)	71.5 _(0.4)	67.0 _(0.0)	76.7 _(0.2)	68.9 _(0.0)
	Single	3	4	64.2 _(0.2)	66.2 _(0.2)	71.3 _(0.2)	64.7 _(0.2)	71.0 _(0.0)	67.0 _(0.0)	76.5 _(0.0)	68.7 _(0.0)
	Multi.	12	1	62.2 _(7.5)	33.8 _(8.3)	52.2 _(3.1)	52.8 _(4.0)	50.8 _(1.2)	55.8 _(4.3)	55.3 _(7.2)	51.9 _(2.2)
	Multi.	6	2	60.0 _(5.5)	33.7 _(8.4)	50.8 _(1.2)	52.2 _(2.4)	50.2 _(0.2)	54.3 _(2.5)	55.0 _(6.7)	50.9 _(1.8)
	Multi.	3	4	58.7 _(4.9)	33.7 _(8.4)	50.8 _(1.2)	51.3 _(1.9)	50.0 _(0.0)	54.3 _(2.5)	55.3 _(7.2)	50.6 _(2.0)

M Broader Impacts

Our findings suggest that existing transformer-based language models possess the ability to learn and adapt to context by internally fine-tuning a complex model *even during inference*. Consequently, although users are unable to directly modify deployed models, these models may still undergo dynamic updates while processing a context left-to-right, resulting in previously unseen behavior by the time the model reaches the end of the context. This has significant implications for the field of model alignment. It is challenging to impose restrictions on a model that can perform such dynamics updates internally, so malicious content can influence the output of deployed models.

Alternatively, we recognize the potential benefits of pre-training constructed models that integrate explicit fine-tuning mechanisms. By embedding the functionalities typically achieved through explicit fine-tuning, such as detecting malicious content and intent within the models themselves, the need for external modules can be mitigated. Pre-training the constructed model may offer a self-contained solution for ensuring safe and responsible language processing without relying on external dependencies.