

TOWARDS LEARNING HIGH-PRECISION LEAST SQUARES ALGORITHMS WITH SEQUENCE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper investigates whether sequence models can learn to perform numerical algorithms, e.g. gradient descent, on the fundamental problem of least squares. Our goal is to inherit two properties of standard algorithms from numerical analysis: (1) *machine precision*, i.e. we want to obtain solutions that are accurate to near floating point error, and (2) *numerical generality*, i.e. we want them to apply broadly across problem instances. We find that prior approaches using Transformers fail to meet these criteria, and identify limitations present in existing architectures and training procedures. First, we show that softmax Transformers struggle to perform high-precision multiplications, which prevents them from precisely learning numerical algorithms. Second, we identify an alternate class of architectures, comprised entirely of polynomials, that can efficiently represent high-precision gradient descent iterates. Finally, we investigate precision bottlenecks during training and address them via a high-precision training recipe that reduces stochastic gradient noise. Our recipe enables us to train two polynomial architectures, gated convolutions and linear attention, to perform gradient descent iterates on least squares problems. For the first time, we demonstrate the ability to train to near *machine precision*. Applied iteratively, our models obtain $100,000\times$ lower MSE than standard Transformers trained end-to-end and they incur a $10,000\times$ smaller generalization gap on out-of-distribution problems. We make progress towards end-to-end learning of numerical algorithms for least squares.

1 INTRODUCTION

Least squares is the workhorse of modern numerics: it is well understood theoretically (Boyd & Vandenberghe, 2004; Trefethen & Bau, 2022) and has important downstream applications in science and engineering, including solving regression problems and differential equations (Orszag, 1972; Trefethen, 2000). Thus, least squares has gained interest as a natural testbed for investigating how well ML models can learn to implement algorithms (Garg et al., 2022; Von Oswald et al., 2023).

A surge of recent work suggests that Transformers (Vaswani et al., 2017) can learn to solve least squares using *optimization algorithms* like gradient descent and Newton’s method (Akyürek et al., 2022; Fu et al., 2023; Ahn et al., 2024; Bai et al., 2024; Zhang et al., 2023b). These arguments rest on two observations: (1) simplified Transformer architectures (e.g. non-causal linear attention) can exactly implement such algorithms; (2) standard (softmax attention) Transformers learn solutions with similar properties (e.g. convergence rates) as iterative algorithms. Crucially, these works focus on *statistical* least squares: they evaluate Transformer solutions in underdetermined/noisy settings and compare to Bayes-optimal estimators. However, scientific applications like climate or fluids modeling require *numerically precise* solutions to least squares, e.g. to accurately model turbulence or to maintain stable temporal rollouts (Frisch, 1995; Wilcox, 2006). Prior works do not engage with the issue of high precision, so it is still unclear how well Transformers can solve least squares from this perspective.

In this work, we thus study whether existing approaches can solve *numerical* least squares. Specifically, numerical analysis requires that algorithms exhibit (1) *machine precision*, i.e. they should obtain solutions that are accurate to near floating point error, and (2) *numerical generality*, i.e. they are computational procedures that should apply broadly across problem instances. (See Section 2.1 for details.) Since traditional least squares algorithms (e.g. gradient descent and conjugate gradients)

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

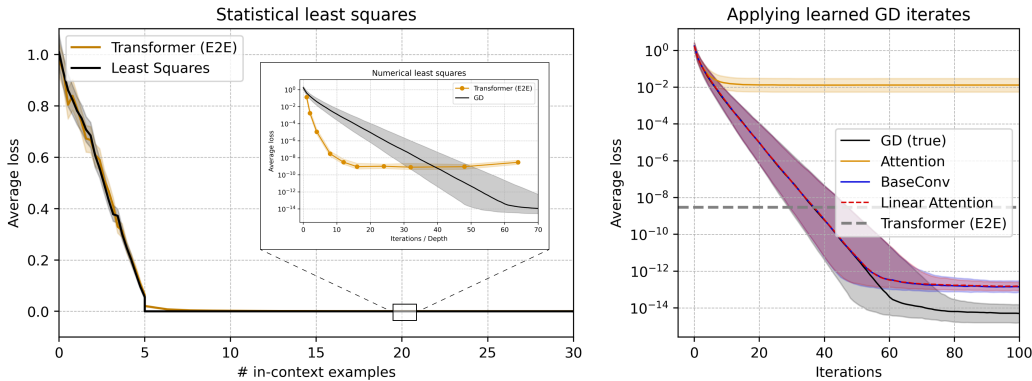


Figure 1: Prior work focuses on *statistical* least squares: Transformers approximate Bayes-optimal estimators (left, adapted from Garg et al. (2022)). In this work, we focus on *numerical* least squares: Transformers struggle to obtain precise solutions (inset). Using a high-precision training recipe, we train two polynomial architectures, BASECONV and linear attention, to perform high-precision gradient descent iterates on least squares (right): applied iteratively, they reach $\approx 10^{-13}$ MSE.

provably meet these criteria (Trefethen & Bau, 2022), it is crucial to evaluate machine learning methods against these same standards to determine their ability to learn numerical algorithms.

We focus on learning the gradient descent (GD) algorithm for least squares. Our study has three parts:

- **We benchmark standard Transformers for precision/generality and identify an expressivity gap on least squares.** When we replicate the standard end-to-end training setup for least squares with Transformers, we find that solutions do not exhibit machine precision and numerical generality (Figure 1a, 2). We identify high-precision multiplications as a fundamental challenge for softmax Transformers. Empirically, on a synthetic element-wise multiplication task, we find precision scales poorly with larger Transformers: an 8-layer model trains to an MSE that is still *10 million* times worse than machine epsilon (Figure 3). Theoretically, we argue that a single layer of softmax attention is unable to exactly express element-wise multiplications. Since implementing GD involves high-precision multiplications, this observation suggests standard Transformers are unable to even precisely *express* GD, much less precisely *learn* the algorithm.
- **We identify an alternate architecture class which does not suffer from expressivity problems.** Motivated by the expressivity limitations of softmax attention, we investigate alternate sequence mixer architectures. Prior work notes that non-causal linear attention is able to exactly implement algorithms like GD and Newton’s method (Von Oswald et al., 2023; Giannou et al., 2024) because it consists entirely of *polynomials*. We provide a unified framework to understand existing expressivity results from the lens of arithmetic circuits. In our work, we focus on BASECONV, a gated convolutional architecture, as a case study, since it is provably equivalent to the entire class of polynomial architectures (Arora et al., 2023; 2024). We demonstrate that gated convolutions can express a high-precision GD algorithm ($\approx 10^{-13}$ MSE when implemented in practice, Figure 4).
- **We identify an optimization precision bottleneck and propose a high-precision training recipe.** Although polynomial architectures can precisely express the GD algorithm, we find that standard training procedures struggle to find a solution with sufficiently high precision (10^{-5} MSE, Figure 9). Therefore, towards disentangling precision bottlenecks during training, we first focus on the intermediate task of explicitly learning GD iterates. We identify stochastic gradient noise from minibatching as the main optimization bottleneck, and we find that a simple metric, cosine similarity of minibatch gradients (Liu et al., 2023b), is diagnostic of precision saturation. Towards reducing stochasticity, we propose (1) a learning rate (LR) scheduler that adaptively adjusts LR based on the cosine similarity metric, and (2) to apply EMA over optimizer updates to maintain strong gradient signal. Our high-precision training recipe allows us to train ML architectures to near *machine precision* for the first time. We successfully train two 3-layer models, with BASECONV and linear attention, that learn to perform a single high-precision iteration of GD (Figure 1b). Excitingly, we can also learn multiple GD iterates at once, scaling up to 4 iterations with 10^{-10} MSE.

Overall, our work makes the following contributions: (1) we specify the desiderata of learning numerical algorithms, *machine precision* and *numerical generality*, and we demonstrate that standard Transformers fall short because of expressivity limitations of softmax attention; (2) we provide a unified framework using arithmetic circuits to investigate the expressivity of the class of polynomial architectures; (3) we address additional precision bottlenecks that emerge *during training*, even when using expressive polynomial architectures. Although we do not achieve end-to-end learning of GD, we make significant headway: we propose a high-precision training recipe, which, for the first time, allows us to learn iterates of the GD algorithm to near *machine precision*.

2 LEARNING NUMERICAL ALGORITHMS FOR LEAST SQUARES

In this section, we distinguish between statistical vs. numerical least squares and discuss the two properties we want our models to inherit from numerical algorithms: *machine precision* and *numerical generality*. We then briefly discuss prior work and, in doing so, tease apart two increasingly end-to-end notions of performing algorithms with ML: *expressing an algorithm in-weights* and *learning algorithm iterates*.

2.1 PROBLEM FORMULATION AND RELATED WORK

In this work, our goal is to train a model that solves least squares problems: find $\mathbf{x} \in \mathbb{R}^D$ given $\mathbf{A} \in \mathbb{R}^{N \times D}$ and $\mathbf{b} \in \mathbb{R}^N$ such that $\mathbf{A}\mathbf{x} = \mathbf{b}$. Here, we briefly discuss two different perspectives on least squares: statistical (in the form of *in-context learning*) and numerical.

Statistical least squares. Originally motivated by applications in language modeling, prior works on solving least squares with Transformers typically take a *statistical* perspective. Transformers are trained using an *in-context learning* setup (Garg et al., 2022; Akyürek et al., 2022): problem instances $\mathbf{A}\mathbf{x} = \mathbf{b}$ are sampled from a pre-specified distribution \mathcal{D}_{train} , and the model is trained to minimize mean squared error (MSE) over \mathcal{D}_{train} . Trained models are then evaluated on unseen problem instances, both in and out-of-distribution, and their performance is compared to Bayes-optimal estimators (Garg et al., 2022; Akyürek et al., 2022). We define the in-context least squares training setup in Appendix B and leave a more detailed discussion of related in-context learning work to Appendix A.

Numerical least squares. In this work, we instead take a *numerical* perspective on least squares. A prototypical numerical algorithm for least squares is GD. For a problem instance $\mathbf{A}\mathbf{x} = \mathbf{b}$, we initialize \mathbf{x}_0 , an estimate of \mathbf{x} , and iteratively improve our estimate via

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla \mathcal{L}(\mathbf{x}_i), \quad (1)$$

where $\mathcal{L}(\hat{\mathbf{x}}) := \frac{1}{2} \|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}\|_2^2$ is the squared residual error. GD exhibits two properties of numerical algorithms that we want our models to inherit:

- *Machine precision.* Numerical algorithms provably obtain high-precision solutions. For GD, obtaining higher precision simply requires performing more iterations until convergence to machine precision (i.e. the smallest achievable error with floating-point arithmetic) (see Chapter 11 of Trefethen & Bau (2022)). In this work, we use `float32` throughout, where machine precision is $2^{-23} \approx 1.19 \times 10^{-7}$, so we hope for MSEs around $2^{-46} \approx 1.42 \times 10^{-14}$.
- *Numerical generality.* Although the convergence rate of GD depends on the spectrum of \mathbf{A} (see Chapter 9 of Boyd & Vandenberghe (2004)), the computational procedure comprising GD is general and can be applied broadly to problem instances. This is unlike statistical generalization and notions of in vs. out-of-distribution. In this work, we are interested to study how closely ML models can emulate the numerical generality of algorithms despite training on a data distribution.

2.2 OUTLINE OF THIS WORK

A recent line of work probes the estimators learned by Transformers on in-context least squares, and suggests that Transformers learn to solve least squares by mimicking iterative algorithms like gradient descent and Newton’s method (Von Oswald et al., 2023; Ahn et al., 2024; Fu et al., 2023;

162 [Giannou et al., 2024](#)). These works typically analyze simplified models theoretically and extrapolate
 163 to standard training regimes, backed by empirical observations:

- 164 • **Theoretical results for simplified models**, e.g. *non-causal linear attention* can implement GD
 165 using a specific choice of model weights.
- 166 • **Empirical experiments training standard Transformers**, e.g. decoder-only softmax attention
 167 Transformers trained *end-to-end* on in-context least squares display convergence rates reminiscent
 168 of iterative algorithms.

169 Although prior works suggest that trained Transformers learn to solve least squares with algorithms, it
 170 is still unclear whether statements about learning algorithms in simplified settings transfer to standard
 171 Transformers trained end-to-end. We note two significant gaps between previously-analyzed settings
 172 and standard in-context least squares:

- 173 • **Architectural differences**. Standard Transformers use softmax instead of linear attention, causal
 174 instead of non-causal sequence mixers, and include MLPs and LayerNorms ([Ba et al., 2016](#)).
- 175 • **Optimization**. Even if a model can express a precise and general algorithm, it is unclear whether
 176 the model can learn the algorithm from data.

177 In this work, we tease apart bottlenecks caused by architecture expressivity limitations (Sections 3.3, 4)
 178 and optimization difficulties (Section 5) by investigating two increasingly sophisticated notions of
 179 performing GD for least squares with ML: *expressing GD in-weights* and *learning GD iterates*.

182 3 TRANSFORMERS DO NOT LEARN NUMERICAL ALGORITHMS IN-CONTEXT

183 In this section, we evaluate standard Transformers, trained end-to-end, on the criteria of *machine*
 184 *precision* and *numerical generality*. Surprisingly, we demonstrate that existing approaches fail to
 185 exhibit these properties: the precision of Transformer solutions saturates $10^6 \times$ worse than machine
 186 precision (Section 3.1), and their performance further degrades as problem instances deviate from the
 187 model’s training distribution (Section 3.2). These results suggest that Transformers are not learning
 188 proper algorithms as numerical analysis defines them.

189 Towards identifying expressivity bottlenecks, we identify three linear algebra primitives that comprise
 190 standard algorithms including GD and Newton’s method (Section 3.3). We find empirically that
 191 Transformers struggle to implement high-precision multiplication, and theoretically we argue that
 192 softmax attention faces an expressivity gap when trying to exactly express multiplications.

195 3.1 TRANSFORMERS STRUGGLE TO REACH MACHINE PRECISION

196 Recent work ([Von Oswald et al., 2023](#); [Ahn et al., 2024](#); [Fu et al., 2023](#); [Giannou et al., 2024](#))
 197 studying in-context least squares suggests that Transformers learn to mimic iterative algorithms like
 198 GD and Newton’s method. Note that if Transformers are able to implement iterative algorithms,
 199 the depth of the model should correspond to the number of iterations performed. We thus focus on
 200 the simplest case of fully determined least squares problems with fixed size design matrices and
 201 investigate whether precision improves as we scale to larger and deeper models.

202 In Figure 1b, following prior work ([Ahn et al., 2024](#)), we fix the size of $\mathbf{A} \in \mathbb{R}^{20 \times 5}$ and train
 203 Transformers end-to-end on least squares, scaling up to $L = 64$ layers. We compare their precision
 204 to the convergence rate of the full-batch gradient descent algorithm on least squares. For more details
 205 about the training setup, refer to Appendix B.3.1.

206 At first, Transformer precision scaling exceeds the convergence rate of gradient descent: this finding
 207 mirrors similar results reported by [Fu et al. \(2023\)](#), who suggest Transformers may instead be learning
 208 higher-order algorithms like Newton’s method. However, we further observe that the precision
 209 gains for Transformers *rapidly diminish*, such that we observe very little difference in precision
 210 between $L = 32$ and $L = 64$ layers. The deepest Transformer models we are able to train achieve
 211 an MSE around 10^{-8} . In contrast, gradient descent converges linearly to machine precision, almost
 212 $1,000,000 \times$ better precision. The diminishing returns of the Transformer precision scaling imply
 213 that Transformers are not learning standard numerical algorithms like GD.

3.2 TRANSFORMERS DO NOT EXHIBIT THE GENERALITY OF GRADIENT DESCENT

We further investigate whether Transformers learn solutions to least squares that exhibit numerical generality. Recall that models are trained on a predefined distribution of least squares problems, \mathcal{D}_{train} . If Transformers learn to solve least squares using a standard numerical algorithm like GD, then we expect the performance of the model should be robust to out-of-distribution inputs.

For GD specifically, the convergence criterion ($0 < \eta < 2/\sigma_{max}^2$) depends on σ_{max} , the maximum singular value of \mathbf{A} , and the optimal rate of convergence depends on the condition number of \mathbf{A} , $\kappa = \sigma_{max}/\sigma_{min}$ (Boyd & Vandenberghe, 2004). Thus we specify our training distribution \mathcal{D}_{train} over least squares problems ($\mathbf{A} \in \mathbb{R}^{20 \times 5}$, $\mathbf{b} = \mathbf{A}\mathbf{x} \in \mathbb{R}^{20}$) as follows. First, as in prior work (Garg et al., 2022), we sample the entries of \mathbf{A} and \mathbf{x} i.i.d from a standard Gaussian $N(0, 1)$. We then shift and rescale the singular values of \mathbf{A} so that $\sigma_{max} = \kappa = 5$. After training a 12-layer Transformer model on the in-context objective, we evaluate our model on out-of-distribution regression targets \mathbf{b} .

We define $\mathcal{D}_{OOD}^b(\sigma)$ by sampling each entry of \mathbf{x} i.i.d. from $N(0, \sigma)$ and computing $\mathbf{b} = \mathbf{A}\mathbf{x}$. Although the distribution of \mathbf{b} 's and \mathbf{x} 's changes with σ , because the spectra of the \mathbf{A} 's is consistent, we know that GD with fixed choice of η will provably converge to high precision.

We find that compared to GD, the Transformer solutions are brittle to unseen regression target distributions. Simply scaling the inputs by a factor of $10\times$, the MSE of the trained Transformer degrades by a factor of 10^8 . In contrast, GD is robust: a fixed number of GD iterations consistently converges to the same order of magnitude of precision (Figure 2). The brittleness of the learned Transformer solution compared to GD again suggests that Transformers are not performing standard numerical algorithms.

3.3 IDENTIFYING AN EXPRESSIVITY GAP WITH STANDARD TRANSFORMERS

Toward understanding the limitations of the Transformer architecture, we start with GD and Newton’s method, two algorithms used to solve least squares, and look into primitives that comprise them.

Linear algebra primitives. We observe that GD and Newton’s method can be expressed as compositions of three simple linear algebra operations: sequence-wise read/write (READ), affine transformations (LINEAR), and element-wise multiplications (MULTIPLY). For input $\mathbf{u} \in \mathbb{R}^{N \times D}$:

$$\text{READ}(i, j, a, b)(\mathbf{u}) = \begin{cases} \mathbf{u}[k, a:b] & k \neq j \\ \mathbf{u}[i, a:b] & k = j \end{cases}$$

$$\text{LINEAR}(\mathbf{H})(\mathbf{u}) = \mathbf{u}\mathbf{H}, \quad \text{where } \mathbf{H} : \mathbb{R}^D \rightarrow \mathbb{R}^{d_{out}} \text{ is linear,}$$

$$\text{MULTIPLY}(a, b, d_{out})(\mathbf{u}) = \mathbf{u}[:, a:a+d_{out}] \odot \mathbf{u}[:, b:b+d_{out}]$$

In Appendix D.2, we define these primitives formally and describe how GD and Newton’s method iterates can each be expressed as a composition of these primitives. Intuitively, READ is required to transfer information across the sequence dimension, LINEAR to transfer information across the hidden dimension, and MULTIPLY to compute high-degree interaction terms (like dot products or element-wise squaring).

Empirical analysis: standard Transformers struggle with multiplication. We train Transformers on synthetic formulations of these tasks to investigate how precision scales with model size. Details about our training setups are in Appendix B.3.2.

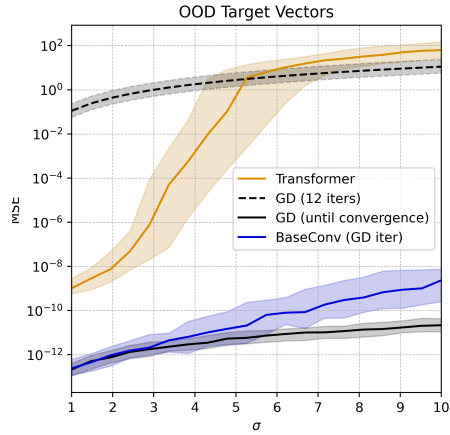


Figure 2: Transformers generalize poorly to out-of-distribution regression targets. In contrast, using our training recipe, we train a BASECONV model to perform high-precision GD iterates. Applied iteratively, our BASECONV model incurs $10,000\times$ less generalization error on out-of-distribution target vectors than the Transformer.

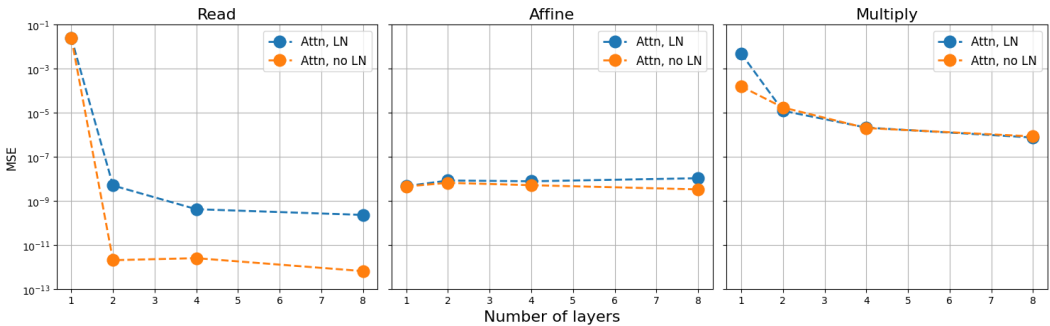


Figure 3: Precision vs. Transformer depth, with and without LayerNorms (LN), on synthetic tasks. While shallow Transformers are able to learn the READ and LINEAR tasks to high precision ($< 10^{-8}$ with 2-layer models), precision on the MULTIPLY task scales poorly with depth (only 10^{-6} with 8-layer models).

In Figure 3, we show that even 2-layer Transformers are able to achieve 10^{-8} MSE on the READ and LINEAR tasks. However, we find that Transformers struggle with the MULTIPLY task: precision scales poorly with model depth, such that an 8-layer Transformer is only able to achieve 10^{-6} MSE. In Appendix C.1, we further show that precision on the MULTIPLY task also scales poorly with increased attention dimension, number of attention heads, and MLP upscaling factor.

Theoretical analysis: softmax attention struggles to exactly express multiplication. In Appendix D.2.4, we provide a proof that a single layer of softmax attention cannot exactly express the simple element-wise squaring function $\text{SQUARE}(\mathbf{u})[i, j] = \mathbf{u}[i, j]^2$ (intuitively, because softmax cannot implement polynomials). Crucially, we note that element-wise squaring is a special case of element-wise multiply, so softmax attention cannot exactly implement MULTIPLY either:

Theorem 3.1 (Informal statement of Theorem D.31 and Corollary D.32). *One-layer single-headed (causal) softmax attention cannot exactly represent SQUARE and MULTIPLY for all possible inputs.*

Since precisely implementing numerical algorithms like GD hinges on performing high-precision multiplications, this result suggests that the standard Transformer architecture struggles to precisely implement these algorithms because of a fundamental expressivity gap.

We mention briefly that these findings do not conflict with prior results (Yun et al., 2020b) proving universal approximation theorems for Transformers, because they typically require parameter count to scale exponentially with dimension: see Appendix A.

4 ALTERNATE ARCHITECTURES CLOSE THE EXPRESSIVITY GAP

Motivated by the finding that softmax attention struggles to precisely express multiplications, we next investigate alternate sequence mixer architectures. We are inspired by prior results (Von Oswald et al., 2023; Giannou et al., 2024) that show non-causal linear attention is able to exactly implement algorithms like GD and Newton’s method. Thus, we focus on the class of *polynomial architectures*, i.e. sequence mixers comprised entirely of polynomial operations, in order to explicitly bake in multiplications. In this section, we present a unified framework that integrates previous findings through the perspective of arithmetic circuits. Specifically, we focus on BASECONV, a gated convolutional model that combines element-wise multiplications (gating) with long convolutions. We work with BASECONV for two reasons:

- Recent work (Arora et al., 2023; 2024) has shown that BASECONV is equivalent to general arithmetic circuits, including all polynomial architectures. Thus, existing results with other polynomial architectures, e.g. linear attention, transfer directly to BASECONV.
- Empirically, gated convolutional models have been shown to perform comparably to attention-based architectures on tasks like language, audio, and DNA modeling (Arora et al., 2024; Nguyen et al., 2024; Zhang et al., 2023a).

We emphasize that although we find gated convolutions are convenient to work with theoretically and empirically, we believe that other sequence mixer architectures may also be able to alleviate the expressivity issues we highlight in Section 3.3. In particular, we show promising empirical results for non-causal linear attention in Section 5.2.

4.1 GATED CONVOLUTIONS ARE EQUIVALENT TO ARITHMETIC CIRCUITS

BASECONV definition. In this work, we focus on a variant of the BASECONV operator from Arora et al. (2023). Given an input $u \in \mathbb{R}^{N \times D}$, $\text{BASECONV}(u)$ is defined as:

$$\underbrace{((uW_{gate} + b_{gate}))}_{\text{Linear Projection}} \odot \underbrace{(h * (uW_{in} + b_{in}) + b_{conv})}_{\text{Convolution}} W_{out} + b_{out} \tag{2}$$

where the layer is parameterized by learnable filters $h \in \mathbb{R}^{N \times D}$, linear projections $W_{in}, W_{gate}, W_{out} \in \mathbb{R}^{D \times D}$, and bias matrices $b_{conv}, b_{in}, b_{gate}, b_{out} \in \mathbb{R}^{N \times D}$. Here, \odot represents the Hadamard product, and convolution of two matrices is computed as convolution of the corresponding columns.

BASECONVs can exactly express linear algebra primitives. In Appendix D.2.1, we provide explicit constructions of single-layer BASECONV models that exactly implement the READ, LINEAR, and MULTIPLY primitives from Section 3.3.

We note that this result is stronger than prior BASECONV expressivity results (e.g. Theorem H.21 from Arora et al. (2023)), which imply a poly-log-factor increase in parameters (specifically layers) translating from arbitrary arithmetic circuits. Here, we show by construction that these specific primitives, and any circuits that are compositions of them, incur only a constant factor loss.

BASECONVs can perfectly recover linear algebra primitives from data. In Appendix D.7, for SQUARE and LINEAR, we further show the following under mild assumptions, which our input distribution satisfies (see details in Assumptions D.45, D.46, D.55, D.61):

Theorem 4.1 (Informal statement of Theorems D.59, D.62). *BASECONV perfectly recovers SQUARE and LINEAR when it achieves zero population gradient w.r.t. MSE loss.*

We note that although results of the form “exact solution implies zero population gradient” exist in the literature (Ahn et al., 2024; Mahankali et al., 2023), to the best of our knowledge, we are the first to show the *converse* (“zero population gradient implies recovery of exact solution”) for sequence model architectures. In Appendix C.1, we show that BASECONV models can learn the READ, LINEAR, and MULTIPLY primitives to high precision in practice (Figure 6).

BASECONVs are universal approximators. Finally, we show in Appendix D.4 that BASECONV can efficiently approximate *smooth functions* by implementing polynomials:

Theorem 4.2 (Informal statement of Theorem D.39). *Given a k -times differentiable function $\bar{f} : [-1, 1] \rightarrow \mathbb{R}$, define $f : [-1, 1]^{N \times D} \rightarrow \mathbb{R}^{N \times D}$, which applies \bar{f} element-wise to all inputs. Then $\forall \epsilon > 0$, there exists a BASECONV model approximates f to within error ϵ , with $O\left(k \sqrt{\frac{L}{\epsilon}}\right) + k$ depth and $O(ND)$ parameters, where $\|f^{(k)}\|_{\infty} \leq L$.*

We additionally prove a universal approximation theorem for general smooth multivariate functions in Appendix D.4 (Theorem D.44).

4.2 BASECONV CAN PRECISELY EXPRESS GRADIENT DESCENT FOR LEAST SQUARES

We now focus on the gradient descent algorithm for least squares. Explicitly, given a least squares problem instance $Ax = b$ and an initial iterate x_0 , a single iteration of gradient descent computes

$$x_1 := x_0 - \eta \nabla \mathcal{L}(x_0), \text{ where } \nabla_x \mathcal{L} = A^T(Ax - b). \tag{3}$$

We provide two explicit $O(1)$ -layer weight constructions to express a GD iterate using BASECONV in Appendix D.3.1. One requires a $O(D)$ state size using a *non-causal* model (i.e. each entry can access

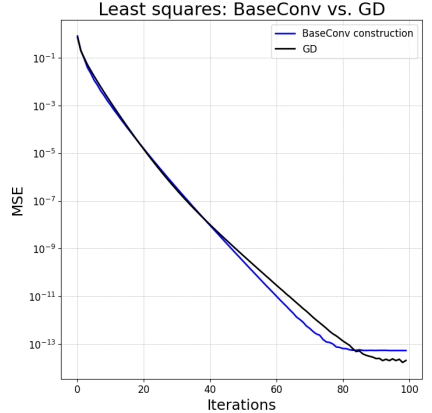


Figure 4: BASECONV can express high-precision gradient descent: our implementation of the weight construction reaches 10^{-13} MSE in practice.

any other entry of the sequence) and one requires a $O(D^2)$ state size using a *causal* model (i.e. entries cannot access later entries of the sequence). In Appendix D.3.2, we prove that both constructions are asymptotically optimal with respect to state size.

In Figure 4, we implement our non-causal weight construction into a deep BASECONV model as a proof of concept. We confirm that gated convolutions can empirically implement high-precision gradient descent – notably, roundoff errors due to machine precision do not significantly accumulate in practice, despite scaling up to a depth-100 BASECONV model.

5 TOWARDS TRAINING MODELS TO MACHINE PRECISION

Although BASECONVs are expressive enough to solve least squares precisely, we find that simply swapping out softmax attention with BASECONV and training end-to-end is insufficient for high precision: our BASECONV models perform as poorly as standard Transformers (Figure 9). This suggests that additional precision bottlenecks are present *during high-precision training*. In this section, we thus investigate what it takes to train polynomial architectures to machine precision.

Recent works (Rodionov & Prokhorenkova, 2023; 2024) on algorithm learning find that intermediate supervision is crucial for learning long computation trajectories. We hypothesize that end-to-end least squares faces a similar challenge. Thus, to study high-precision optimization, we first investigate a simplified setting: learning to perform *explicit GD updates* for least squares.

Using this task as a benchmark, we identify a fundamental bottleneck in high-precision regimes, *gradient variance from minibatching*, and we identify a metric based on cosine similarity of successive gradients that is diagnostic of precision saturation during training. We then propose a high-precision training recipe, which for the first time allows us to train ML models to near *machine precision*. Using our training recipe, we learn to perform explicit GD updates to 10^{-13} average MSE (Figure 1b), and we can also learn up to 4 iterates of GD at once with an MSE of 10^{-10} (Table 7).

Simplifying the training setup. We first define a sequence of *k-th iterate* tasks, where the goal is to explicitly produce the *k*-th iterate of GD given a least squares problem instance (\mathbf{A}, \mathbf{b}) , an initial iterate \mathbf{x}_0 , and a step size η :

$$\{(\mathbf{a}_1, b_1), \dots, (\mathbf{a}_N, b_N), \mathbf{x}_0\} \rightarrow \mathbf{x}_k, \text{ where } \mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla \mathcal{L}(\mathbf{x}_i), i \in [k-1]. \quad (4)$$

We then define the *explicit gradient* task, where the goal is to produce the GD update vector:

$$\{(\mathbf{a}_1, b_1), \dots, (\mathbf{a}_N, b_N), \mathbf{x}_0\} \rightarrow \nabla \mathcal{L}(\mathbf{x}_0). \quad (5)$$

Note that (up to a residual connection), the explicit gradient task is equivalent to 1-step GD, and standard in-context least squares is equivalent to taking $k \rightarrow \infty$. Thus, the explicit gradient task is a natural simplification of standard in-context least squares, and the *k*-th iterate task allows us to smoothly interpolate between the two extremes of difficulty. Refer to Appendix B for more details.

5.1 TOWARDS A HIGH-PRECISION TRAINING RECIPE

Our theoretical results in Section 4 imply that a 3-layer BASECONV is expressive enough to solve the explicit gradient task, so we use training a 3-layer BASECONV on this task as our benchmark for studying the challenges of high-precision learning.

Precision saturates with standard training procedures. Motivated by prior work (Garg et al., 2022; Von Oswald et al., 2023; Ahn et al., 2024), we start by investigating two basic optimization procedures: Adam with constant learning rate (LR) and with exponentially decaying LR.

In Appendix C.2 (Figure 10), we sweep initial LR and LR step rate across 2-3 orders of magnitude for constant and decaying LR schedules. We find:

- *Precision saturation occurs with both constant and decaying LR schedules.* After a number of training iterations, the average loss saturates and is unable to improve. We note that this occurs even while gradients magnitudes and LR are non-zero.
- *Slower-decaying LR schedules perform better but require exponentially more training iterations.* In Figure 8, we further analyze this phenomenon in the simpler case of 1-layer Transformers/BASECONVs on the MULTIPLY synthetic. We observe a power-law relation between precision

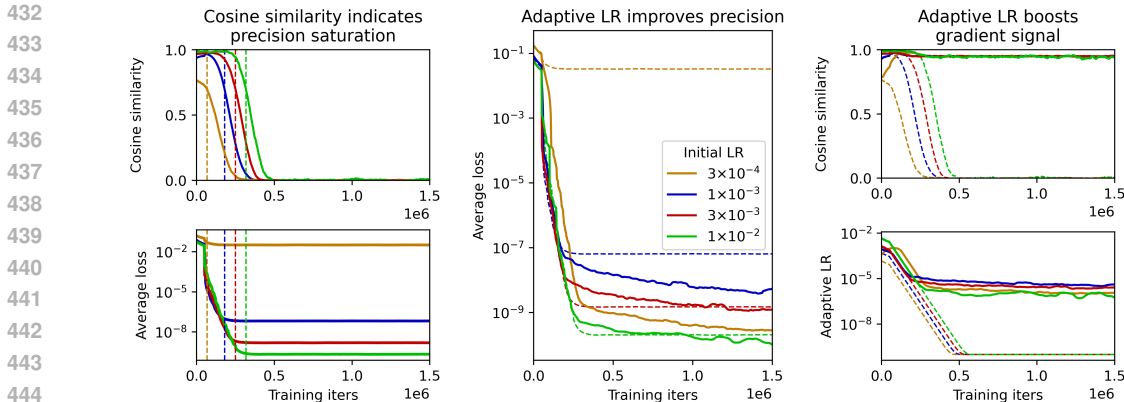


Figure 5: Gradient metric is predictive of precision saturation (left). We propose a simple adaptive LR scheduler that alleviates precision saturation (middle). Adaptive LR effectively boosts gradient signal during training (right).

and number of training iterations as we sweep steprate; although it may be possible to train to high precision in theory, this approach seems infeasible in practice.

- *With aggressively-decaying LR schedules, higher initial LR is better.* For a fixed scheduler step rate, increasing initial LR leads to significant improvements in final MSE, e.g. in Figure 10, an improvement of 1000× simply by increasing initial LR from 10⁻³ to 10⁻². Choosing a LR that is too large causes training instability, so in practice we set it to the *largest value that trains stably*.

Our analysis suggests that Adam with an exponentially decreasing LR scheduler gets us only part of the way to a machine precision training recipe. We next address the issue of precision saturation.

Stochastic gradients bottleneck precision. We identify *minibatch gradient variance* as the main source of precision saturation. Although our goal is to minimize the expected loss over problem instances from \mathcal{D}_{train} , in practice we minimize over finite minibatch samples instead. Minibatch training is the standard in ML, but interestingly we find that the variance in minibatch gradients can dominate the population gradient signal in high-precision regimes, causing the loss to stagnate.

To demonstrate this, we define a simple metric to assess the strength of the gradient signal during training. At a given training step, we take the current model weights, sample n different minibatches of least squares problems, and compute the minibatch model gradients $\{g_1, \dots, g_n\}$. We then compute the *average cosine similarity* between all pairs, as in Liu et al. (2023b):

$$\sigma_g := \frac{2}{n(n-1)} \sum_{i \neq j} \frac{g_i^T g_j}{\|g_i\|_2 \|g_j\|_2} \tag{6}$$

We observe that this cosine similarity metric is predictive of precision saturation across MSE scales and optimizer hyperparameters (Figure 5).

An adaptive LR scheduler boosts gradient signal beyond precision saturation. We thus propose an *adaptive* LR scheduler based on the gradient variance. Our scheduler is motivated by two intuitions:

- Whenever the cosine similarity metric is *high*, gradient signal is strong. In order to refine the highest-precision bits of the model weights, we need to slowly *decrease* the LR.
- Whenever the cosine similarity metric is *low*, the model weights are stuck in a local region of the loss landscape. To allow the model to escape this region, we need to *increase* the LR.

The basic scheduler we use in this work simply decreases the LR exponentially while the metric is above a threshold σ_{th} and increases the LR instead if the metric is below σ_{th} . In Figure 5, we show that this simple approach alleviates the loss saturation phenomenon: we see a boost in population loss across our LR settings, and we observe the models consistently improve as we continue training. We note that proper choice of LR hyperparameters is still crucial for *efficient* convergence to machine precision – we leave speeding up the convergence rate via better adaptive schedulers to future work.

Exponential Moving Average (EMA) over optimizer updates. Finally, motivated by our observation that gradient variance bottlenecks precision and inspired by recent works (Lee et al., 2024; Pagliardini et al., 2024), we apply an additional EMA over Adam’s update vectors to help smooth out minibatch noise. Empirically, we find this boosts the final MSE by as much as $100,000\times$ on the explicit gradient task: see Appendix C.2 (Figure 11).

Our training recipe for efficient high-precision convergence thus involves two techniques: (1) an adaptive LR scheduler that exponentially increases or decays LR according to the cosine similarity metric; and (2) applying EMA over optimizer updates.

5.2 LEARNING HIGH-PRECISION GRADIENT DESCENT WITH POLYNOMIAL ARCHITECTURES

Using our training recipe, we successfully train two 3-layer models with polynomial architectures, BASECONV and non-causal linear attention, on the explicit gradient task. For the first time, we are able to train to near *machine precision*: we achieve an average loss of 10^{-13} MSE.

In Figure 1b, we slot our trained models into the standard GD algorithm, using their predictions in place of the true least squares gradients $\nabla\mathcal{L}$. Specifically, for a least squares problem $\mathbf{Ax} = \mathbf{b}$ and initial iterate \mathbf{x}_0 , we repeatedly compute $\mathbf{x}_{i+1} := \mathbf{x}_i - \eta\Delta_i$, where $\Delta_i := T_\theta(\mathbf{A}, \mathbf{b}, \mathbf{x}_i)$ is the prediction of the model. We iteratively apply the model until convergence to a fixed point \mathbf{x}_∞ .

We find that both our models achieve high precision. In this setting, we reach an average MSE of 10^{-12} (Figure 1, right): this is $100,000\times$ better MSE than the biggest Transformers we are able to train end-to-end. Moreover, our BASECONV model exhibits better numerical generality than the Transformer, incurring a $10,000\times$ smaller generalization gap on problems outside its training distribution (Figure 2). Interestingly, we find that our linear attention model exhibits markedly worse generality: its out-of-distribution performance nearly matches the Transformer’s, and the model iterates eventually diverge: see Figure 13.

Learning k -iterates of GD for larger k . We find that our training recipe also allows us to learn up to $k = 4$ iterates of GD at once with 10^{-10} MSE: see Table 7 and Figure 14 for results. We are not able to stably train deeper models without reintroducing non-polynomial normalization techniques like LayerNorms, which causes precision bottlenecks. For small k , we observe that LayerNorms worsen precision by over $1,000\times$. See Appendix C.3 for details.

Experiments with in-context ODE solving. Finally, towards high-precision ML for more realistic tasks, we provide preliminary results on in-context ODE solving. We find that our proposed techniques outperform standard Transformers by up to $1,000,000\times$ in MSE (up to $\approx 10^{-10}$ with iterative BASECONVs vs. $\approx 10^{-4}$ with 12-layer Transformers). See Appendix C.4 for details.

6 DISCUSSION AND LIMITATIONS

In this work, we investigate learning to solve least squares from a numerical perspective. We find that Transformers fail to learn solutions that exhibit the properties of machine precision and numerical generality. Disentangling effects from the model architecture and optimizer, we find that standard design choices perform surprisingly poorly from the lens of numerics. We identify expressivity limitations with softmax attention, and find surprisingly that even MLPs and LayerNorms significantly affect precision (up to $1,000,000\times$ worse MSE on the explicit gradients task). On the optimization front, we find stochastic gradient noise from minibatch training becomes a precision bottleneck in high-precision regimes. We propose an adaptive LR scheduler that alleviates this issue on a simplified task, but we suspect that this issue remains a fundamental challenge on harder problems. Crucially, although we make progress toward learning to solve numerical least squares end-to-end, our techniques struggle to maintain stable and precise training with deep networks.

We note that the *numerical* criteria we consider in this work represent a fundamentally different type of learning and generalization from *statistical* notions that are prevalent in ML. We believe these numerical perspectives may be relevant to the wider scientific ML community. For example, existing approaches to solving PDEs have shown promise but are known to be brittle outside their training distributions (Wang & Lai, 2023; Rathore et al., 2024). This inhibits their usefulness in high-impact applications like climate or fluids modeling, where high precision and robustness are crucial. We believe learning to implement precise numerical algorithms directly from data is an exciting prospect that has the potential to unlock new capabilities across science and engineering.

540 REPRODUCIBILITY STATEMENT

541
542 We will release the code and configuration files necessary to reproduce our experiments upon
543 acceptance of this paper. In this work, all experiments are done using synthetic data and tasks. All
544 experiments were conducted using PyTorch on NVIDIA A100/H100 GPUs. Detailed hyperparameters
545 (learning rate, batch size, and optimizer settings) and proofs of all theoretical claims are provided in
546 the supplementary materials.

547
548 REFERENCES

- 549 Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement
550 preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing*
551 *Systems*, 36, 2024.
- 552
553 Kabir Ahuja, Madhur Panwar, and Navin Goyal. In-context learning through the bayesian prism.
554 *arXiv preprint arXiv:2306.04891*, 2023.
- 555
556 Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algo-
557 rithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*,
558 2022.
- 559
560 Ekin Akyürek, Bailin Wang, Yoon Kim, and Jacob Andreas. In-context language learning: Architec-
561 tures and algorithms. *arXiv preprint arXiv:2401.12973*, 2024.
- 562
563 Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri Rudra,
564 and Christopher Ré. Zoology: Measuring and Improving Recall in Efficient Language Models,
565 2023.
- 566
567 Simran Arora, Sabri Eyuboglu, Michael Zhang, Aman Timalsina, Silas Alberti, Dylan Zinsley,
568 James Zou, Atri Rudra, and Christopher Ré. Simple linear attention language models balance the
569 recall-throughput tradeoff, 2024.
- 570
571 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. URL
572 <https://arxiv.org/abs/1607.06450>.
- 573
574 Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians:
575 Provable in-context learning with in-context algorithm selection. *Advances in neural information*
576 *processing systems*, 36, 2024.
- 577
578 Stephen P Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge university press, 2004.
- 579
580 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
581 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
582 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 583
584 Wuyang Chen, Jialin Song, Pu Ren, Shashank Subramanian, Dmitriy Morozov, and Michael W
585 Mahoney. Data-efficient operator learning via unsupervised pretraining and in-context learning.
586 *arXiv preprint arXiv:2402.15734*, 2024.
- 587
588 Xiang Cheng, Yuxin Chen, and Suvrit Sra. Transformers implement functional gradient descent to
589 learn non-linear functions in context, 2024. URL <https://arxiv.org/abs/2312.06528>.
- 590
591 David Chiang, Peter Cholak, and Anand Pillay. Tighter bounds on the expressivity of transformer
592 encoders. In *International Conference on Machine Learning*, pp. 5544–5562. PMLR, 2023.
- 593
594 Liam Collins, Advait Parulekar, Aryan Mokhtari, Sujay Sanghavi, and Sanjay Shakkottai. In-
595 context learning with transformers: Softmax attention adapts to function lipschitzness, 2024. URL
596 <https://arxiv.org/abs/2402.11639>.
- 597
598 D. Jackson. *The theory of approximation*. Amer. Math. Soc. Colloq. Publ., vol. 11, Amer. Math. Soc,
599 Providence, R. I., 1930.

- 594 Tri Dao, Nimit S Sohoni, Albert Gu, Matthew Eichhorn, Amit Blonder, Megan Leszczynski, Atri
595 Rudra, and Christopher Ré. Kaleidoscope: An efficient, learnable representation for all structured
596 linear maps. *arXiv preprint arXiv:2012.14966*, 2020.
- 597
598 Ishita Dasgupta, Andrew K Lampinen, Stephanie CY Chan, Antonia Creswell, Dharshan Kumaran,
599 James L McClelland, and Felix Hill. Language models show human-like content effects on
600 reasoning. *arXiv preprint arXiv:2207.07051*, 2022.
- 601 Uriel Frisch. *Turbulence: the legacy of AN Kolmogorov*. Cambridge university press, 1995.
- 602
603 Daniel Y Fu, Tri Dao, Khaled K Saab, Armin W Thomas, Atri Rudra, and Christopher Ré.
604 Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint*
605 *arXiv:2212.14052*, 2022.
- 606
607 Deqing Fu, Tian-Qi Chen, Robin Jia, and Vatsal Sharan. Transformers learn higher-order optimization
608 methods for in-context learning: A study with linear models. *arXiv preprint arXiv:2310.17086*,
609 2023.
- 610
611 Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn
612 in-context? a case study of simple function classes. *Advances in Neural Information Processing*
613 *Systems*, 35:30583–30598, 2022.
- 614
615 Angeliki Giannou, Shashank Rajput, Jy-yong Sohn, Kangwook Lee, Jason D Lee, and Dimitris
616 Papailiopoulos. Looped transformers as programmable computers. In *International Conference on*
617 *Machine Learning*, pp. 11398–11442. PMLR, 2023.
- 618
619 Angeliki Giannou, Liu Yang, Tianhao Wang, Dimitris Papailiopoulos, and Jason D Lee. How well
620 can transformers emulate in-context newton’s method? *arXiv preprint arXiv:2403.03183*, 2024.
- 621
622 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
623 *preprint arXiv:2312.00752*, 2023.
- 624
625 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured
626 state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- 627
628 Michael T Heideman and C Sidney Burrus. *Multiplicative complexity, convolution, and the DFT*.
629 Springer, 1988.
- 630
631 Maximilian Herde, Bogdan Raonić, Tobias Rohner, Roger Käppli, Roberto Molinaro, Emmanuel
632 de Bézenac, and Siddhartha Mishra. Poseidon: Efficient foundation models for pdes, 2024. URL
633 <https://arxiv.org/abs/2405.19101>.
- 634
635 Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers. *arXiv preprint*
636 *arXiv:2310.05249*, 2023.
- 637
638 Jaerin Lee, Bong Gyun Kang, Kihoon Kim, and Kyoung Mu Lee. Grokfast: Accelerated grokking by
639 amplifying slow gradients, 2024. URL <https://arxiv.org/abs/2405.20233>.
- 640
641 Jerry Weihong Liu, N Benjamin Erichson, Kush Bhatia, Michael W Mahoney, and Christopher Re.
642 Does in-context operator learning generalize to domain-shifted settings? In *The Symbiosis of Deep*
643 *Learning and Differential Equations III*, 2023a.
- 644
645 Zhuang Liu, Zhiqiu Xu, Joseph Jin, Zhiqiang Shen, and Trevor Darrell. Dropout reduces underfitting,
646 2023b. URL <https://arxiv.org/abs/2303.01500>.
- 647
648 Arvind Mahankali, Tatsunori B Hashimoto, and Tengyu Ma. One step of gradient descent is
649 provably the optimal in-context learner with one layer of linear self-attention. *arXiv preprint*
650 *arXiv:2307.03576*, 2023.
- 651
652 Nick McGreivy and Ammar Hakim. Weak baselines and reporting biases lead to overoptimism in
653 machine learning for fluid-related partial differential equations. *Nature Machine Intelligence*, 6
654 (10):1256–1269, September 2024. ISSN 2522-5839. doi: 10.1038/s42256-024-00897-5. URL
655 <http://dx.doi.org/10.1038/s42256-024-00897-5>.

- 648 William Merrill and Ashish Sabharwal. A logic for expressing log-precision transformers. *Advances*
649 *in Neural Information Processing Systems*, 36, 2024.
- 650 William Merrill and Ashish Sabharwals. The parallelism tradeoff: Limitations of log-precision
651 transformers. *Transactions of the Association for Computational Linguistics*, 11:531–545, 2023.
652 doi: 10.1162/tacl_a_00562. URL <https://aclanthology.org/2023.tacl-1.31>.
- 653 Eric J. Michaud, Ziming Liu, and Max Tegmark. Precision machine learning. *Entropy*, 25(1):175,
654 January 2023. ISSN 1099-4300. doi: 10.3390/e25010175. URL <http://dx.doi.org/10.3390/e25010175>.
- 655 Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures
656 for grokking via mechanistic interpretability, 2023. URL <https://arxiv.org/abs/2301.05217>.
- 657 Eric Nguyen, Michael Poli, Matthew G. Durrant, Armin W. Thomas, Brian Kang, Jeremy Sul-
658 livan, Madelena Y. Ng, Ashley Lewis, Aman Patel, Aaron Lou, Stefano Ermon, Stephen A.
659 Baccus, Tina Hernandez-Boussard, Christopher Ré, Patrick D. Hsu, and Brian L. Hie. Se-
660 quence modeling and design from molecular to genome scale with evo. *bioRxiv*, 2024.
661 doi: 10.1101/2024.02.27.582234. URL <https://www.biorxiv.org/content/early/2024/02/27/2024.02.27.582234>.
- 662 Steven A Orszag. Comparison of pseudospectral and spectral approximation. *Studies in Applied*
663 *Mathematics*, 51(3):253–259, 1972.
- 664 Matteo Pagliardini, Pierre Ablin, and David Grangier. The ademamix optimizer: Better, faster, older,
665 2024. URL <https://arxiv.org/abs/2409.03137>.
- 666 Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Huanqi Cao, Xin
667 Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, et al. Rwkv: Reinventing rnns for the
668 transformer era. *arXiv preprint arXiv:2305.13048*, 2023.
- 669 Peter Bürgisser and Michael Clausen and M. Amin Shokrollah. *Algebraic Complexity Theory*.
670 Springer, 1997.
- 671 Tobias Von Petersdorff. Polynomial approximation and interpolation. 2015. Numerical Analysis Class
672 Notes. <https://www.math.umd.edu/~petersd/666/amsc666notes02.pdf>.
- 673 W. Pleśniak. Multivariate jackson inequality. *Journal of Computational and Applied Math-*
674 *ematics*, 233(3):815–820, 2009. ISSN 0377-0427. doi: <https://doi.org/10.1016/j.cam.2009.02.095>. URL <https://www.sciencedirect.com/science/article/pii/S0377042709001307>. 9th OPSFA Conference.
- 675 Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua
676 Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional
677 language models. In *International Conference on Machine Learning*, pp. 28043–28078. PMLR,
678 2023.
- 679 Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: General-
680 ization beyond overfitting on small algorithmic datasets, 2022. URL <https://arxiv.org/abs/2201.02177>.
- 681 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
682 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 683 Pratik Rathore, Weimu Lei, Zachary Frangella, Lu Lu, and Madeleine Udell. Challenges in training
684 pinns: A loss landscape perspective, 2024. URL <https://arxiv.org/abs/2402.01868>.
- 685 Allan Raventós, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the
686 emergence of non-bayesian in-context learning for regression. *Advances in Neural Information*
687 *Processing Systems*, 36, 2024.
- 688 Gleb Rodionov and Liudmila Prokhorenkova. Neural algorithmic reasoning without intermediate
689 supervision, 2023. URL <https://arxiv.org/abs/2306.13411>.

- 702 Gleb Rodionov and Liudmila Prokhorenkova. Discrete neural algorithmic reasoning, 2024. URL
703 <https://arxiv.org/abs/2402.11628>.
704
- 705 Günther Schulz. Iterative berechnung der reziproken matrix. *ZAMM-Journal of Applied Mathematics
706 and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik*, 13(1):57–59, 1933.
- 707 Smoothness. Smoothness — Wikipedia, the free encyclopedia, 2006. [https://en.wikipedia.
708 org/wiki/Smoothness](https://en.wikipedia.org/wiki/Smoothness).
- 709 Gilbert Strang. *Linear algebra and its applications*. 2012.
710
- 711 Lloyd N Trefethen. *Spectral methods in MATLAB*. SIAM, 2000.
- 712 Lloyd N. Trefethen and David Bau. *Numerical Linear Algebra, Twenty-fifth Anniversary
713 Edition*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2022. doi:
714 10.1137/1.9781611977165. URL [https://epubs.siam.org/doi/abs/10.1137/1.
715 9781611977165](https://epubs.siam.org/doi/abs/10.1137/1.9781611977165).
- 716 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
717 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing
718 systems*, 30, 2017.
719
- 720 Petar Veličković, Adrià Puigdomènech Badia, David Budden, Razvan Pascanu, Andrea Banino,
721 Misha Dashevskiy, Raia Hadsell, and Charles Blundell. The clrs algorithmic reasoning benchmark,
722 2022. URL <https://arxiv.org/abs/2205.15659>.
- 723 Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,
724 Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In
725 *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.
726
- 727 Yongji Wang and Ching-Yao Lai. Multi-stage neural networks: Function approximator of machine
728 precision, 2023. URL <https://arxiv.org/abs/2307.08934>.
- 729 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
730 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.
731 *arXiv preprint arXiv:2206.07682*, 2022.
- 732 Sanford Weisberg. *Applied linear regression*, volume 528. John Wiley & Sons, 2005.
733
- 734 D.C. Wilcox. *Turbulence Modeling for CFD*. Number v. 1 in Turbulence Modeling for CFD. DCW
735 Industries, 2006. ISBN 9781928729082. URL [https://books.google.com/books?id=
736 tFNNPgAACAAJ](https://books.google.com/books?id=tFNNPgAACAAJ).
- 737 Steve Yadlowsky, Lyric Doshi, and Nilesh Tripuraneni. Pretraining data mixtures enable narrow
738 model selection capabilities in transformer models. *arXiv preprint arXiv:2311.00871*, 2023.
- 739 Liu Yang, Siting Liu, Tingwei Meng, and Stanley J. Osher. In-context operator learning with data
740 prompts for differential equation problems. *Proceedings of the National Academy of Sciences*,
741 120(39), September 2023a. ISSN 1091-6490. doi: 10.1073/pnas.2310142120. URL [http:
742 //dx.doi.org/10.1073/pnas.2310142120](http://dx.doi.org/10.1073/pnas.2310142120).
- 743 Liu Yang, Siting Liu, Tingwei Meng, and Stanley J Osher. In-context operator learning for differential
744 equation problems. *arXiv preprint arXiv:2304.07993*, 2023b.
745
- 746 Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J. Reddi, and Sanjiv Kumar. Are
747 transformers universal approximators of sequence-to-sequence functions?, 2020a.
- 748 Chulhee Yun, Yin-Wen Chang, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv
749 Kumar. O (n) connections are expressive enough: Universal approximability of sparse transformers.
750 *Advances in Neural Information Processing Systems*, 33:13783–13794, 2020b.
- 751 Michael Zhang, Khaled K. Saab, Michael Poli, Tri Dao, Karan Goel, and Christopher Ré. Effectively
752 modeling time series with simple discrete state spaces, 2023a. URL [https://arxiv.org/
753 abs/2303.09489](https://arxiv.org/abs/2303.09489).
- 754 Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context.
755 *arXiv preprint arXiv:2306.09927*, 2023b.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

APPENDIX

The appendix is organized as follows:

- Appendix [A](#) provides a more detailed overview of related work.
- Appendix [B](#) provides details about our experimental setup.
- Appendix [C](#) provides additional experiments and ablation studies.
- Appendix [D](#) provides details about our main theoretical results.

810 A EXTENDED BACKGROUND

811 A.1 LEAST SQUARES

812 Least squares, $\mathbf{Ax} = \mathbf{b}$, is well-understood theoretically, and we know of simple numerical algorithms
813 for solving least squares to high precision (Weisberg, 2005; Boyd & Vandenberghe, 2004). We focus
814 on two algorithms: gradient descent and Newton’s method.

815 **Gradient descent** Given a guess for \mathbf{x}^* , we minimize the least squares loss

$$816 \mathcal{L}(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{a}_i^T \mathbf{x} - b_i)^2 \quad (7)$$

817 via gradient descent on \mathbf{x} :

$$818 \nabla_{\mathbf{x}} \mathcal{L}_N = \sum_{i=1}^N (\mathbf{x}^T \mathbf{a}_i - b_i) \mathbf{a}_i \quad (8)$$

$$819 \mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla \mathcal{L}_N(\mathbf{x}_t) \quad (9)$$

820 **Ordinary Least Squares and Newton’s method** In the noiseless, full determined regime, the
821 Bayes-optimal estimator is ordinary least squares (OLS) (Weisberg, 2005):

$$822 \mathbf{x}^{OLS} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}, \quad (10)$$

823 where

$$824 \mathbf{A} = \begin{pmatrix} \leftarrow \mathbf{a}_1 \rightarrow \\ \vdots \\ \leftarrow \mathbf{a}_N \rightarrow \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} b_1 \\ \vdots \\ b_N \end{pmatrix} \quad (11)$$

825 Note that this estimator requires a matrix inverse, which is expensive to compute exactly. An
826 alternative is to use Newton’s method to approximate the matrix inverse term (Schulz, 1933). To
827 estimate $(\mathbf{A}^T \mathbf{A})^{-1}$, we can perform the following iterative algorithm:

$$828 \mathbf{M}_{t+1} = \mathbf{M}_t (2\mathbf{I} - (\mathbf{A}^T \mathbf{A}) \mathbf{M}_t) \quad (12)$$

829 where \mathbf{M}_t converges to $(\mathbf{A}^T \mathbf{A})^{-1}$.

830 A.2 RELATED WORK

831 In this section, we detail prior work on in-context learning, Transformer expressivity, gated convolu-
832 tional architectures, and algorithm learning.

833 **In-context learning.** The capability of Transformers to perform in-context learning on language
834 and pattern matching tasks has been well-documented (Brown et al., 2020; Dasgupta et al., 2022; Wei
835 et al., 2022). More recently, a flurry of work has investigated in-context learning for regression-style
836 tasks. Garg et al. (2022) first formulated the mathematical framework to analyze the estimators
837 Transformers implement in-context, focusing on linear regression and other least squares problems. A
838 number of works further observed empirically that Transformers seem to approximate Bayes-optimal
839 estimators on distributional problems. For example, based on the task distribution, the performance of
840 in-context Transformers mimics optimally-tuned LASSO on sparse linear regression, ridge regression
841 on noisy dense linear regression, and Bayes-optimal priors for task mixtures (Akyürek et al., 2024;
842 Raventós et al., 2024; Yadlowsky et al., 2023; Ahuja et al., 2023; Bai et al., 2024). Beyond standard
843 least squares problems, other works have investigated the ability of Transformers to in-context solve
844 broader problems of scientific interest like differential equations (Yang et al., 2023b; Chen et al.,
845 2024; Liu et al., 2023a).

846 Towards explaining these observations, recent works have focused on understanding the expressivity
847 and optimization landscapes of Transformer variants (typically non-causal linear attention) on linear
848 regression. Linear attention has been shown to be expressive enough to implement numerical
849 algorithms for solving linear regression, including gradient descent (Akyürek et al., 2022; Von Oswald

et al., 2023) and Newton’s method (Fu et al., 2023; Giannou et al., 2024). Recent work (Ahn et al., 2024; Mahankali et al., 2023; Zhang et al., 2023b) has also begun to investigate the optimization dynamics for linear attention on least squares. Finally, we highlight that recent work (Bai et al., 2024; Huang et al., 2023; Collins et al., 2024; Cheng et al., 2024) makes progress on theoretically understanding non-linear attention, e.g. with softmax or ReLU activations.

Unlike prior work, we investigate the capabilities of standard Transformers, focusing on exploring their capability to perform *high-precision* optimization algorithms. Noting a gap between empirical performance and theoretical claims regarding in-context least squares as gradient descent, we further investigate alternative architectures to softmax attention.

Expressivity and approximation ability of Transformers. Although Transformers were initially designed for discrete tasks like language modeling, recent works have investigated the ability of the Transformer architecture to express general *continuous-valued* sequence-to-sequence maps. We briefly mention three classes of prior work:

- **Constructive arguments.** We highlight Giannou et al. (2023), which proposes a looped-Transformer weight construction that implements a basic mathematical instruction set. Using compositions of these instructions, the authors demonstrate that Transformers are expressive enough to implement numerical algorithms, including matrix inversion and SGD on linear models.
- **Universal approximation results.** Several works, such as Yun et al. (2020a;b), provide bounds on the number of parameters and layers required to approximate smooth sequence-to-sequence functions to arbitrary precision using Transformers. However, these results typically require parameters to scale exponentially with respect to problem size, which quickly becomes impractical in practice.
- **Complexity theory results.** Recent works (Chiang et al., 2023; Merrill & Sabharwal, 2023; Merrill & Sabharwal, 2024) prove that *log-precision* Transformers lie in TC^0 , a limited complexity class of circuits.

Gated convolutions. Gated convolutional models are a class of architectures that serve as an efficient alternative to attention. These models, consisting of gating (element-wise multiplication) and long convolutions (filter size equal to sequence length), stem from earlier work (Gu et al., 2021) inspired by the signal processing literature. In this work we focus on the BASECONV model from Arora et al. (2023), but a recent surge of interest in efficient attention replacements has led to a flood of gated convolutional architectures (Poli et al., 2023; Peng et al., 2023; Gu & Dao, 2023).

Recent architectural innovations within the class of gated convolutional models have been largely motivated by language modeling tasks (Fu et al., 2022; Arora et al., 2023). Unlike these prior works, which focus on matching attention’s performance on *discrete* tasks, we observe that the connection between gated convolutions and arithmetic circuits implies they are able to exactly express a range of important numerical algorithms for *continuous-valued* tasks. We further investigate their ability to learn these algorithms in-context.

Algorithm learning. We mention two lines of work related to learning algorithms using ML:

- **Grokking.** Several works (Power et al., 2022; Nanda et al., 2023; Lee et al., 2024) have observed the ability of Transformers to learn to perfectly perform small discrete algorithmic tasks, e.g. modular arithmetic.
- **Neural Algorithmic Reasoning.** Recent work (Rodionov & Prokhorenkova, 2023; 2024) investigates the ability of graph neural networks to learn fundamental algorithms like breadth-first search (Veličković et al., 2022).

Crucially, we note that these previous works focus on learning *discrete* algorithmic tasks, which Transformers excel at. As far as we know, we are the first to investigate whether Transformers are able to learn *numerical* algorithms, which rely on addressing key challenges with high-precision floating-point arithmetic.

918 **Precision and scientific ML.** The importance and difficulty of high-precision ML for scientific
919 settings is well-established: although the scientific ML community has made exciting progress in
920 recent years, numerical methods are still known to outperform existing ML methods in precision
921 even on simple PDE benchmarks (McGreivy & Hakim, 2024). Despite this, we are aware of only a
922 few works which directly focus on investigating high precision for ML. We highlight (Michaud et al.,
923 2023; Wang & Lai, 2023), which focus on small MLPs for regression tasks and propose alternate
924 training recipes.

925 As far as we are aware, we are the first to investigate and isolate effects of model architectures and
926 optimizers on precision in a controlled setting: in-context least squares. We find that typical training
927 recipes for sequence models (e.g. softmax attention, Adam, and standard LR schedulers) encounter
928 surprising precision barriers when applied to numerical tasks.

929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

972 B EXPERIMENTAL SETUP

973 Here, we provide additional details about our experimental setup.

974 B.1 MODEL ARCHITECTURE

975 We base our Transformer and BASECONV models off the GPT2 family (Radford et al., 2019). Unless
976 otherwise specified, we use the following default settings for Transformers:

977 Config	978 Setting
979 Embedding size	64
980 Number of layers	12
981 Number of heads	1
982 MLPs	True
983 MLP hidden size	4× embedding size
984 MLP activation	ReLU
985 LayerNorms	True
986 Input dim	5
987 Sequence length	20

988 Table 1: Standard Transformer architecture details.

989 and the following settings for BASECONVs:

990 Config	991 Setting
992 Embedding size	64
993 Number of layers	3
994 MLPs	False
995 LayerNorms	False
996 Input dim	5
997 Sequence length	20

998 Table 2: BASECONV architecture details.

999 Finally, we describe the settings we use for our linear attention experiment (Figure 2):

1000 Config	1001 Setting
1002 Embedding size	256
1003 Number of layers	3
1004 Number of heads	16
1005 MLPs	False
1006 LayerNorms	False
1007 Input dim	5
1008 Sequence length	20

1009 Table 3: Linear attention architecture details.

1010 B.2 OPTIMIZER

1011 We describe two sets of optimizer settings we use throughout this work.

1012 The first, representative of standard training procedures, is inspired by prior in-context learning
1013 setups (Garg et al., 2022; Von Oswald et al., 2023).

1014 The second, our training recipe, is for our high-precision experiments, where we find a more
1015 aggressive learning rate scheduler is essential. Note we use the adaptive learning rate scheduler and
1016 EMA described in Section 5.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

Config	Setting
Batch size	256
Optimizer	Adam
Learning rate	10^{-3}
Scheduler	StepLR
Training iterations	10^6
Step rate	10^4
Decay rate	0.9

Table 4: Standard optimizer settings.

Config	Setting
Batch size	1024
Optimizer	Adam
Learning rate	10^{-2}
Scheduler	AdaptiveLR
Training iterations	2.5×10^6
Step rate	3×10^3
Decay rate	0.9
EMA decay	0.98
EMA lambda	2

Table 5: High-precision training recipe settings for BASECONV.

Finally, we describe the optimization settings we used for high-precision linear attention, which we found needed a slightly different learning rate scheduler.

Config	Setting
Batch size	1024
Optimizer	Adam
Learning rate	10^{-2}
Scheduler	AdaptiveLR
Training iterations	2.5×10^6
Step rate	3×10^3
Decay rate	0.9
EMA decay	0.98
EMA lambda	2

Table 6: High-precision training recipe settings for linear attention.

1080 B.3 TASKS

1081

1082 Each of our in-context learning tasks can be viewed as a sequence-to-sequence map

1083

1084

$$\mathcal{M} : \mathbb{R}^{N_{in} \times D_{in}} \rightarrow \mathbb{R}^{N_{out} \times D_{out}}$$

1085

1086 In this subsection, we provide details about task implementations, specifying the input/output formats
1087 for each of the synthetic tasks and in-context least squares variants we implement.

1088

1088 B.3.1 IN-CONTEXT LEAST SQUARES.

1089

1090

1091 We consider $\mathcal{M}_{LS} : \mathbb{R}^{N \times (D+1)} \rightarrow \mathbb{R}^D$, where as above the inputs are formatted as

1091

1092

$$\mathbf{u}_{in} := \begin{bmatrix} \mathbf{a}_1 & \dots & \mathbf{a}_N \\ b_1 & \dots & b_N \end{bmatrix}$$

1093

1094 and the expected output is

1095

$$T_\theta(\mathbf{u}_{in})[:-1, -1:] := \mathbf{x}.$$

1096

1097

1097 B.3.2 PRIMITIVES.

1098

1099

1098 For each of the following linear algebra primitives, we increase the task size, setting $D = 20$ and
1099 $N = 40$.

1100

1101

- READ is defined as $\mathcal{M}_{Read} : \mathbb{R}^{N \times D} \rightarrow \mathbb{R}^{N \times D}$, where the inputs are formatted as

1102

1103

$$\mathbf{u}_{in} \in \mathbb{R}^{N \times D} := [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_N]$$

1104

1105 and the expected outputs are $T_\theta(\mathbf{u}_{in}) \in \mathbb{R}^{N \times D}$ such that

1106

1107

$$T_\theta(\mathbf{u}_{in})[k, :] := \begin{cases} \mathbf{u}_{in}[i, :] & k = j \\ \mathbf{u}_{in}[k, :] & k \neq j \end{cases}$$

1108

1109 for task parameters $i \neq j \in [N]$.

1110

- LINEAR is defined as $\mathcal{M}_{Linear} : \mathbb{R}^{N \times D} \rightarrow \mathbb{R}^{N \times 1}$, where the inputs are formatted as

1111

1112

$$\mathbf{u}_{in} \in \mathbb{R}^{N \times D} := [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_N]$$

1113

1114 and the expected outputs are

1115

$$T_\theta(\mathbf{u}_{in}) := [\mathbf{x}_1^T \mathbf{h} \quad \dots \quad \mathbf{x}_N^T \mathbf{h}]$$

1116

1117 where $\mathbf{h} \in \mathbb{R}^D$ is a task parameter.

1118

1119

- MULTIPLY is defined as $\mathcal{M}_{Multiply} : \mathbb{R}^{N \times D} \rightarrow \mathbb{R}^{N \times D/2}$, where the inputs are formatted
1120 as

1121

1122

$$\mathbf{u}_{in} \in \mathbb{R}^{N \times D} := [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_N]$$

1123

1124 and the expected outputs are

1125

$$T_\theta(\mathbf{u}_{in}) := (\mathbf{x}_1[:, : D/2] \odot \mathbf{x}_1[:, D/2 :]) \quad \dots \quad \mathbf{x}_N[:, : D/2] \odot \mathbf{x}_N[:, D/2 :]).$$

1126

1127

1127 B.3.3 EXPLICIT GRADIENT UPDATES.

1128

1129

1128 In Section 5, we investigate a simple training setting, in which the model is explicitly trained
1129 to predict the gradient of the least squares loss. We proceed to define the task $\mathcal{M}_{gradient} :$
1130 $\mathbb{R}^{(N+1) \times (D^2 + 2D + 1)} \rightarrow \mathbb{R}^D$.

1131

1132 The inputs are formatted as

1133

$$\mathbf{u}_{in} := \begin{bmatrix} \mathbf{a}_1 & \dots & \mathbf{a}_N & \mathbf{x}_0 \\ b_1 & \dots & b_N & 0 \end{bmatrix}.$$

1134

1135 The expected outputs are

$$T_\theta(\mathbf{u}_{in})[-1:, :D] := \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{x}_0).$$

1134 B.3.4 k -TH GRADIENT DESCENT ITERATE.

1135
1136 Finally, toward end-to-end least squares, we investigate a series of increasingly end-to-end tasks in
1137 which the model is explicitly trained to predict the k -th gradient descent iterate. We proceed to define
1138 the task $\mathcal{M}_{iter}^k : \mathbb{R}^{(N+1) \times (D^2+2D+1)} \rightarrow \mathbb{R}^D$.

1139 The inputs are formatted as

$$1140 \mathbf{u}_{in} := \begin{bmatrix} \mathbf{a}_1 & \dots & \mathbf{a}_N & \mathbf{x}_0 \\ b_1 & \dots & b_N & 0 \end{bmatrix}.$$

1143 The expected outputs are

$$1144 T_{\theta}(\mathbf{u}_{in})[-1:, :D] := \mathbf{x}_k.$$

1146 B.4 DATA GENERATION

1147
1148 At each training step, we produce a random training prompt \mathbf{u}_{in} by sampling each variable randomly:
1149 from the isotropic Gaussian distribution $N(\mathbf{0}, \mathbf{I})$ for continuous-valued parameters, and from the
1150 uniform distribution for discrete parameters. Concretely:

- 1151 • For the in-context linear regression tasks, input vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$ are sampled from
1152 $N(\mathbf{0}^D, \mathbf{I}^D)$, and the unknown linear function is determined by \mathbf{w}^* , also drawn from
1153 $N(\mathbf{0}^D, \mathbf{I}^D)$.
- 1154 • For the synthetic tasks READ, LINEAR, MULTIPLY (Section 3.3), *each column* of the inputs
1155 $\mathbf{u}_{in} \in \mathbb{R}^{N \times D}$ is sampled from the isotropic Gaussian distribution $N(\mathbf{0}^D, \mathbf{I}^D)$. The tasks
1156 READ and LINEAR require specifying additional parameters as follows:
1157 – For READ, at each iteration, $i \neq j \in [N]$ are sampled uniformly.
1158 – For LINEAR, at each iteration, the affine transformation \mathbf{h} is sampled from
1159 $N(\mathbf{0}^D, 3\mathbf{I}^D)$.
- 1160 • For the explicit gradient task and the k -th gradient descent iterate task, the random initializa-
1161 tion \mathbf{w}_0 is also drawn from $N(\mathbf{0}^D, \mathbf{I}^D)$.

1163 The model is trained to minimize mean squared error over the distribution of prompts.

1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

C ADDITIONAL EXPERIMENTAL RESULTS

C.1 ABLATIONS: LINEAR ALGEBRA PRIMITIVES

In Figure 6, we train Transformers and BASECONVs, *with* MLPs, with and without LayerNorms (LN), on the READ, LINEAR, and MULTIPLY primitives from Section B.3.2. We vary the model depth $L \in \{1, 2, 4, 8\}$ and investigate how precision scales with number of layers. In these experiments, we use a standard exponentially decaying LR schedule for Adam.

We show that Transformers and BASECONVs both achieve high precision ($< O(10^{-9})$) on the READ and LINEAR tasks. However, the Transformers struggle to implement MULTIPLY to high precision, and performance scales poorly with model depth. We observe that BASECONV without LayerNorm generally performs the best across all three primitives, consistently outperforming BASECONV with LayerNorm by 2-4 orders of magnitude.

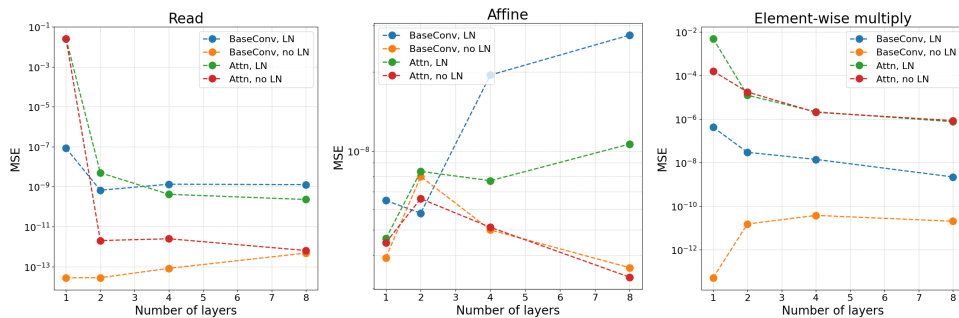


Figure 6: Attention vs. BASECONV, with and without LayerNorms, on synthetic tasks. Precision consistently scales better with depth for BASECONV models than for Transformers. While both models solve READ and LINEAR tasks to at least 10^{-8} MSE, the precision of Transformers scales poorly for the MULTIPLY task.

Focusing on 2-layer Transformers and the MULTIPLY task, we additionally find that precision scales poorly with multiple scaling axes, including hidden dimension, number of heads, and MLP upscaling factor (Figure 7).

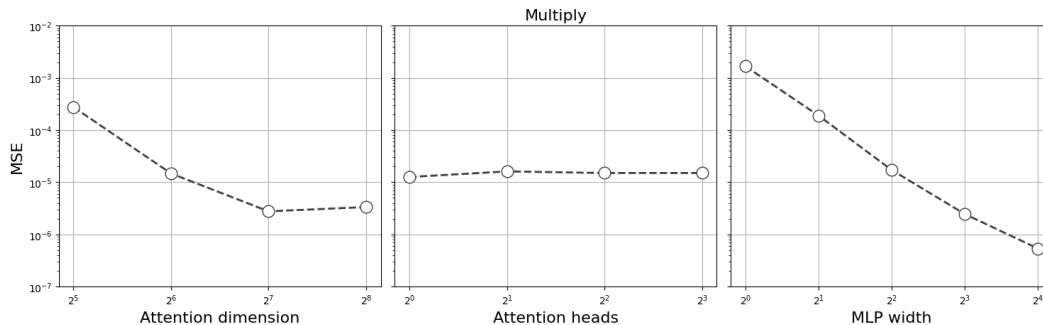


Figure 7: Precision of (2-layer) Transformers on MULTIPLY task scales poorly with attention dimension (left), number of heads (middle), and MLP width (right, where MLP hidden dimension = width \times attention dimension).

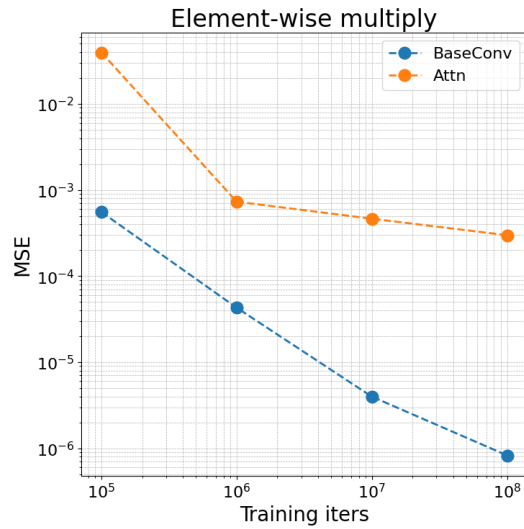


Figure 8: Scaling number of training iterations for 1-layer Transformer vs. BASECONV on the MULTIPLY task. Both models improve precision by 2-3 orders of magnitude as training duration increases by 3 orders of magnitude.

Finally, we investigate the effect of training duration on precision. In Figure 8, we train 1-layer Transformers and BASECONVs, with MLPs and LayerNorms, on the MULTIPLY primitive and vary the number of iterations for which the model is trained. Recall that since new data is sampled at each iteration, we also effectively scale the dataset size proportionally. To keep the learning rates consistent across runs, we scale back the scheduler step size accordingly:

$$\begin{aligned} num_iters &\in \{10^5, 10^6, 10^7, 10^8\} \\ step_size &\in \{10^3, 10^4, 10^5, 10^6\} \end{aligned}$$

We observe a power law, particularly clearly for BASECONV, as we scale from 10^5 to 10^8 iterations. Both models achieve a 2-3 order of magnitude improvement in precision, but this requires also increasing training duration by 3 orders of magnitude.

C.2 ABLATIONS: HIGH-PRECISION OPTIMIZATION

In Figure 9, we try directly training on the end-to-end least squares task, simply replacing softmax attention with BASECONV in the standard Transformer architecture. We find we are unable to reach high precision using this training procedure.

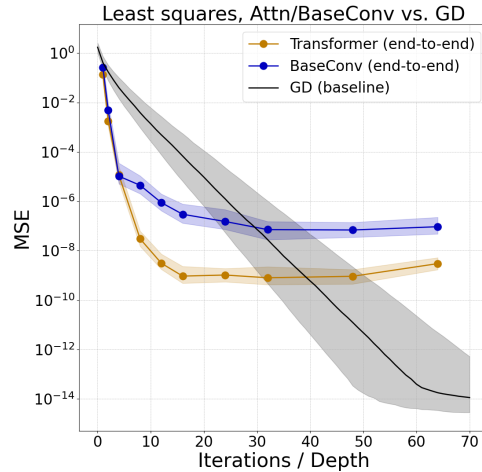


Figure 9: Replacing softmax attention with BASECONV in the standard Transformer architecture and training end-to-end on least squares is not enough to achieve high-precision solutions. BASECONV models trained end-to-end perform as badly as Transformers at small scale, and our largest models perform $100\times$ worse than parameter-matched Transformers.

In Figure 10, we ablate the effects of constant and exponentially decaying LR schedulers with Adam (cutting off training after 10^6 iterations). We find that neither are able to efficiently train to machine precision on the explicit gradients task. For exponentially decaying LR schedule, we find that the LR step rate is a crucial parameter: on the explicit gradient task, a difference of $10,000\times$ between precision saturation thresholds using 1×10^3 vs 3×10^3 for example.

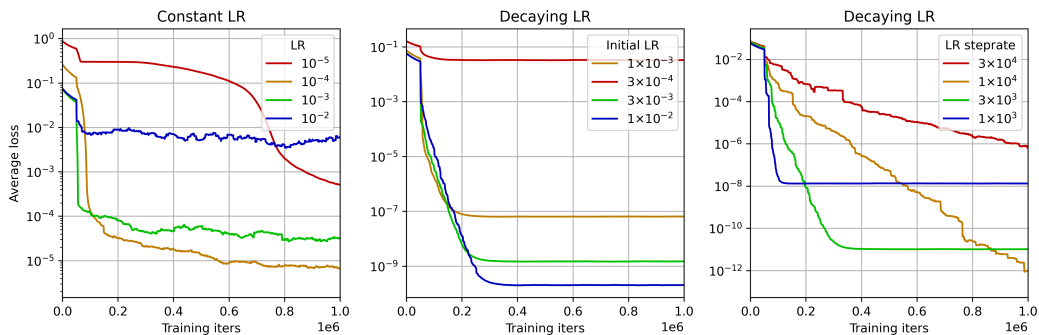


Figure 10: Training with Adam on the explicit gradient task, we ablate LR for constant scheduler (left), initial LR (middle) and LR step rate (right) for decaying scheduler.

In Figure 11, we ablate the effect of applying an EMA over the update vectors from the Adam optimizer. Empirically, we find that this boosts the final MSE by as much as 100,000 \times on the explicit gradient task.

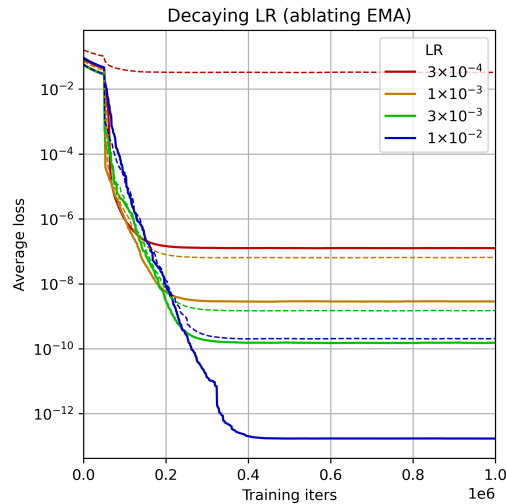


Figure 11: Training on the explicit gradient task, applying EMA over Adam’s update vectors consistently boosts final MSE, by up to 5 orders of magnitude.

In Figure 12, we ablate the effect of restoring the MLPs and LayerNorms to BASECONV models. Surprisingly, we find that even these architectural components worsen the model’s precision: on the explicit gradient task, by a factor of up to 1,000,000 \times MSE. We note that due to training instability with the BASECONV+MLP model, we used a less aggressive LR schedule with initial LR 10^{-3} and LR stepsize 10^4 .

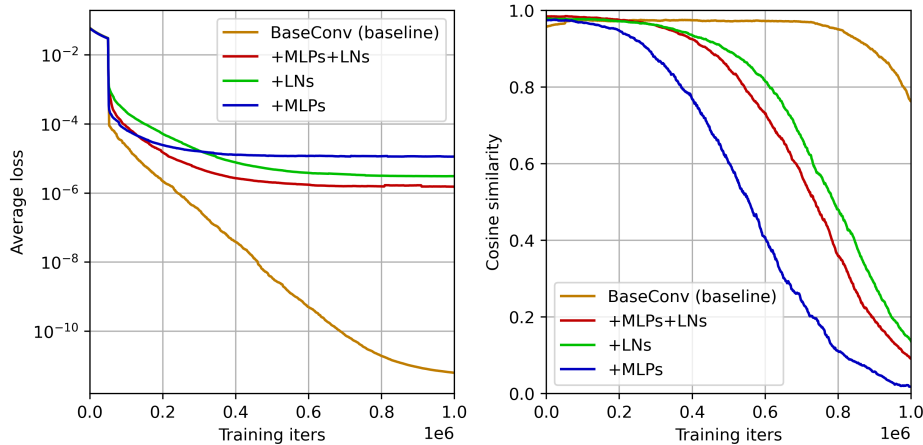


Figure 12: Training on the explicit gradient task, adding MLPs and LNs consistently bottlenecks precision: here, by up to 6 orders of magnitude.

In Figure 13, we evaluate 3-layer BASECONV and linear attention models trained on the explicit gradient task. As in Section 5.2, we apply them iteratively until convergence. We then evaluate on out-of-distribution regression targets, as in Section 3.2.

We surprisingly find that linear attention demonstrates poor numerical generality, despite training to near machine precision on the training distribution. Beyond $\sigma = 4$, the iterations of linear attention diverge.

This result suggests that although different polynomial architectures may equally be able to *express* algorithms, they may *learn* solutions that exhibit vastly different numerical properties.

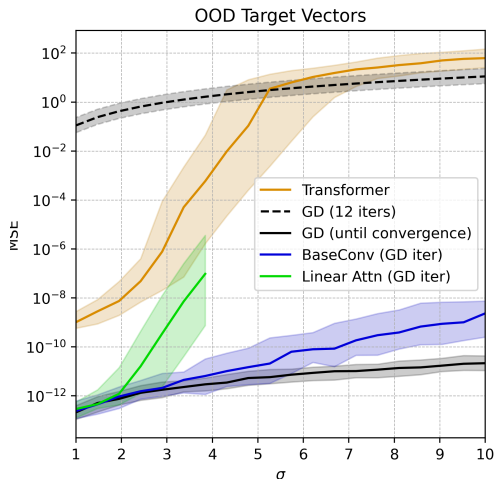


Figure 13: While BASECONV demonstrates improved numerical generality compared to end-to-end trained Transformers, the generalization gap for linear attention is as bad as the Transformer.

C.3 k -TH ITERATE GD

In this section, we investigate how well our proposed techniques can learn the k -th GD iterate tasks as defined in Section 5:

$$\{(\mathbf{a}_1, b_1), \dots, (\mathbf{a}_N, b_N), \mathbf{x}_0\} \rightarrow \mathbf{x}_k, \text{ where } \mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla \mathcal{L}(\mathbf{x}_i), i \in [k - 1]. \quad (13)$$

Recall that $k = 1$ is equivalent to the explicit gradient task, while taking $k \rightarrow \infty$ is equivalent to the standard in-context least squares task. Here, we are interested in understanding how well our techniques extend to larger k , towards learning end-to-end least squares. See Appendix B for a more detailed description of the training setup.

Our theoretical results in Section 4 imply that a $k + 2$ -layer BASECONV is expressive enough to solve the k -th iterate task to machine precision. Thus we train $k + 2$ -layer BASECONV models on the k -th iterate task for $k \geq 1$ using our training recipe.

k	1	2	3	4
MSE	5.0×10^{-13}	2.5×10^{-11}	2.5×10^{-11}	3.1×10^{-10}

Table 7: We can learn up to 4 iterations of GD at once with our current training techniques. Model stability becomes a bottleneck with harder tasks.

Training on the k -iter GD task, we find our training recipe scales to $k = 4$ before training instability occurs. Adding LayerNorms, we are able to train deeper models, but we find MSE worsens by at least $1,000\times$ for small k : see Figure 14.

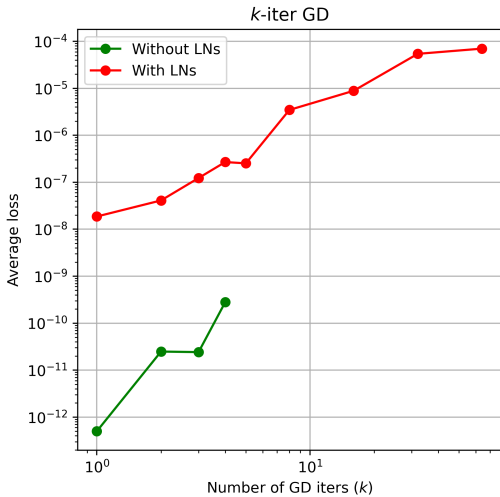


Figure 14: BASECONV with LayerNorms are able to stably scale to deeper models, but LayerNorms present a precision bottleneck: even on small k , MSE degrades by over $1,000\times$.

1512 C.4 IN-CONTEXT ODE SOLVING

1513
 1514 In this section, we demonstrate the generality of our insights on the more practical setting of in-context
 1515 ODE solving. We note that solving differential equations in-context with Transformers is a framework
 1516 that has been explored in recent papers (Yang et al., 2023a; Herde et al., 2024; Liu et al., 2023a),
 1517 and thus represents a natural first step towards extending our techniques to realistic scientific ML
 1518 problems.

1519 **Experimental setup.** We follow the setup from Liu et al. (2023a):

- 1520 • We train on a distribution of 1D ODEs over $t \in [-1, 1]$, defined by

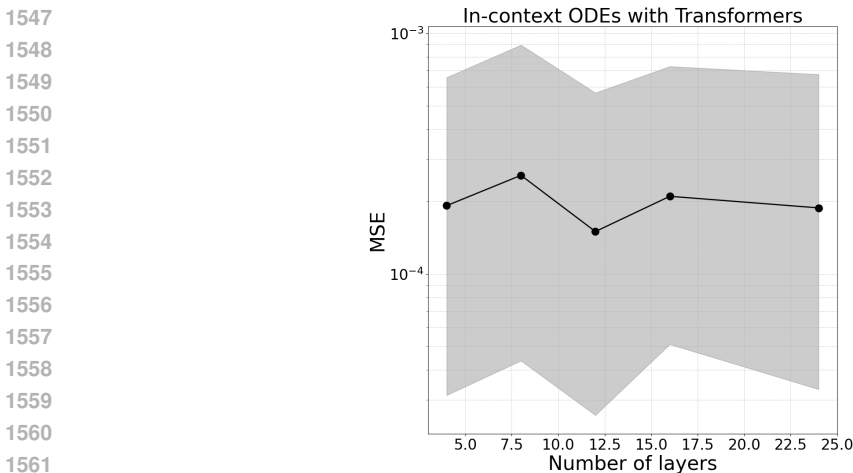
$$1521 \quad u'(t) = \alpha_1 c(t) + \alpha_2 u(t) + \alpha_3. \tag{14}$$

1522
 1523 For each operator, we provide 25 in-context examples of forcing functions, initial conditions,
 1524 and their corresponding solution values at a fixed time $t_{query} \in [-1, 1]$. We then give
 1525 the model a query forcing function and initial condition, and the goal is to predict the
 1526 corresponding solution at t_{query} .

- 1527 • We sample our parameters $\alpha_1 \sim \text{Unif}([0.5, 1.5])$, $\alpha_2 \sim \text{Unif}([-1, 1])$, $\alpha_3 \sim \text{Unif}([-1, 1])$.
- 1528 • Initial conditions are sampled from $u(0) \sim \text{Unif}([-1, 1])$.
- 1529 • Forcing functions $c(t)$ are sampled from a Gaussian process with RBF kernel $K(x, x') =$
 1530 $\exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$, with length-scale parameter $\ell = 1$. We sample each forcing function on
 1531 21 equispaced points over $[-1, 1]$.
- 1532 • ODEs are solve pseudospectrally on $N = 41$ nodes: we find this is sufficient for machine-
 1533 precision solutions with FLOAT32 datatype.

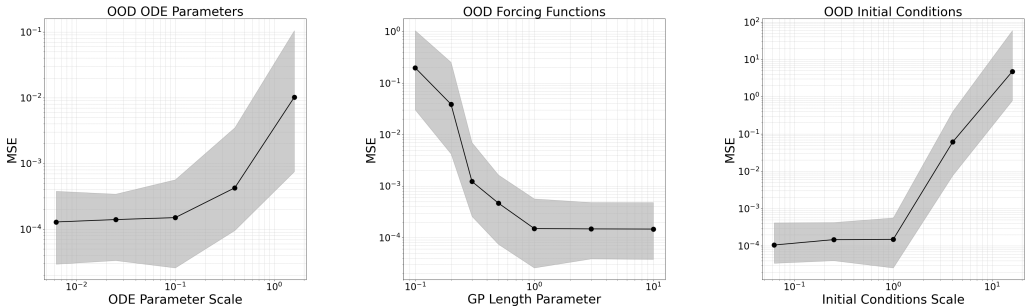
1534
 1535 We find that our observations from least squares transfer to the setting of in-context ODEs:

1536
 1537 **Transformers struggle to learn precise solutions.** We find that a 12-layer, 9M parameter Trans-
 1538 former model only achieves $\approx 10^{-4}$ MSE, almost $10^{10} \times$ worse than the threshold FLOAT32 machine
 1539 epsilon implies. Furthermore, as with least squares, we observe precision saturation with model size.
 1540 In Figure 15, we find that scaling the depth of the model by up to $2 \times$ does not improve precision.
 1541 We further note that precision saturations already seems to occur with 4-layer Transformers. We
 1542 hypothesize that the depth at which precision saturation begins is dependent on the task difficulty.



1543 Figure 15: Transformers fail to learn precise algorithms for solving ODEs in-context. As with
 1544 least squares, precision saturates with deeper models: in our experiments, we observe no significant
 1545 performance boost between 4-layer and 24-layer Transformers.

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576



1577
1578
1579
1580
1581

Figure 16: Transformers fail to learn numerically general solutions: performance is brittle to out-of-distribution ODE parameters (left), forcing function smoothness (middle), and initial condition distribution (right).

1582
1583
1584
1585

Transformers exhibit brittle generalization. We observe that Transformers are not robust to changes to the distributions of ODE parameters, forcing functions, and initial conditions. We describe our experimental setup below, mirroring Liu et al. (2023a):

1586
1587
1588
1589
1590
1591
1592
1593
1594
1595

- **Out-of-distribution ODE parameters.** We parameterize out-of-distribution ODEs via a scale parameter σ_{op} , where $\alpha_1 \sim \text{Unif}([1 - \frac{1}{2}\sigma_{op}, 1 + \frac{1}{2}\sigma_{op}])$ and $\alpha_2, \alpha_3 \sim \text{Unif}([- \sigma_{op}, \sigma_{op}])$. As we increase σ_{op} , we sample from a wider distribution of ODE solution operators, including those with larger operator norms and worse-conditioned design matrices.
- **Out-of-distribution forcing functions.** We vary ℓ , the length parameter of the Gaussian process from which we sample our forcing functions, which effectively controls their smoothness.
- **Out-of-distribution initial conditions.** We sample out-of-distribution initial conditions as $u(0) \sim \text{Unif}([- \sigma_{IC}, \sigma_{IC}])$. As we vary σ , we widen the distribution of the solution values at $t = 0$, which increases the overall magnitudes of the solutions.

1596
1597
1598

We note that in all out-of-distribution experiments, the Transformer’s MSE explodes to near $O(1)$: refer to Figure 16.

1599
1600
1601
1602
1603
1604

Our proposed techniques obtain precise and general solutions. Liu et al. (2023a) shows that in-context ODEs can be reduced to solving least squares problems. Thus, we train a 3-layer BASECONV architecture on the explicit gradient task for the equivalent least squares problem, and apply our model iteratively, as in Section 5. We compare the performance of our iterative model with end-to-end Transformers, least squares solvers, and standard gradient descent applied to the equivalent least squares problem.

1605
1606
1607
1608
1609
1610
1611
1612
1613

We note that our ODEs reduce to least squares problems that are ill-conditioned. In this set of experiments, we find the condition numbers of our design matrices are $O(10^8)$. Since the theoretical convergence rate of gradient descent on least squares depends inversely on the condition number (Boyd & Vandenberghe, 2004), we expect our iterative models and standard gradient descent will require orders of magnitude more iterations than in the least squares problems of Section 5. As such, we limit the number of iterations for our BASECONV model and standard gradient descent to 10,000. Nonetheless, we find that our BASECONV model learns to high enough precision that we are able to maintain the stability of the iterative algorithm for up to $O(10^5)$ steps. In our experiments, we iteratively apply our BASECONV model until convergence to a fixed point and report final MSEs.

1614
1615
1616
1617
1618
1619

We find that BASECONV learns a precise and general algorithm for in-context ODEs:

- **Precision.** In Figure 17, we show that our BASECONV model, applied iteratively, converges to $\approx 10^{-10}$ MSE, $\approx 1,000,000\times$ higher precision than our best Transformers.
- **Generality.** In Figure 18, we find that our BASECONV model exhibits more robust generalization than the Transformer model: in all the out-of-distribution settings we test, our BASECONV model achieves higher precision than Transformers in-distribution. Like above,

we evaluate on out-of-distribution ODE parameters, forcing functions, and initial conditions. In particular, we note that the performance of our BASECONV model almost exactly matches proper gradient descent, even in out-of-distribution settings. Additionally, we find the generalization behavior of our BASECONV model matches the generalization of a proper least squares solver with preconditioning, except for out-of-distribution initial conditions, where we note that the iterative procedure suffers from slow convergence and times out at 10,000 iterations.

We believe these preliminary results show the promise of our techniques towards learning numerical algorithms for more complex tasks, such as solving PDEs, directly from data.

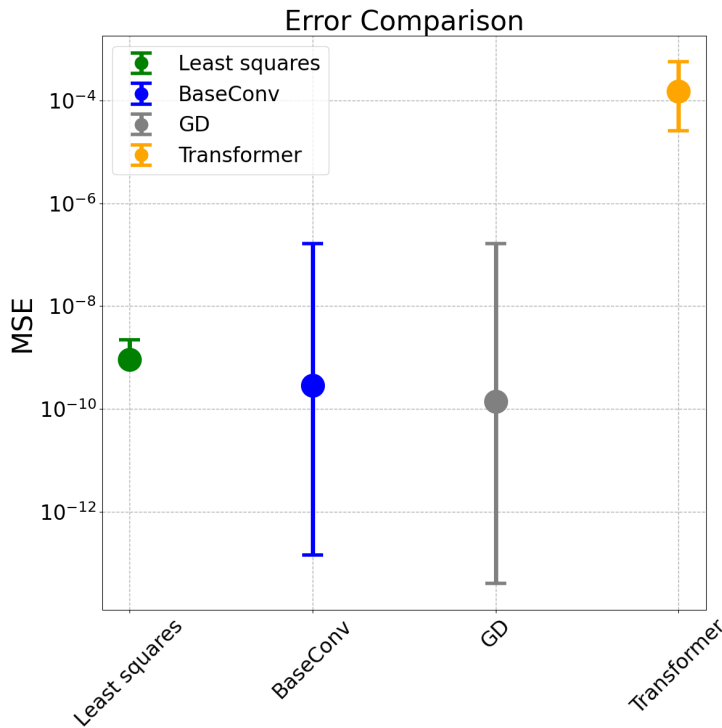


Figure 17: In-distribution error comparison between Transformer, BASECONV, gradient descent, and least squares.

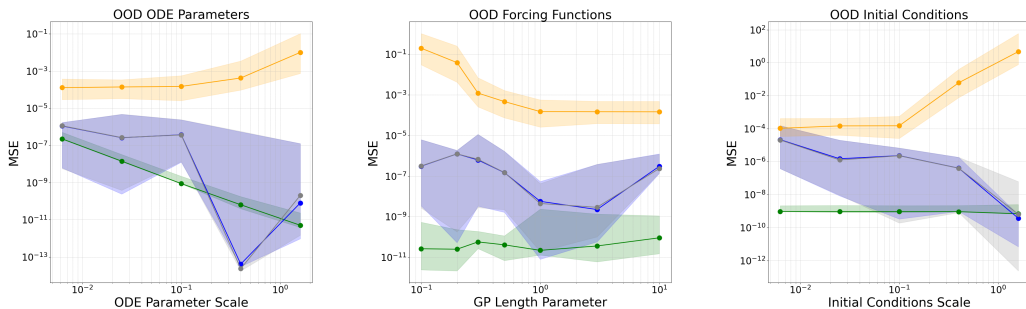


Figure 18: Out-of-distribution error comparison between Transformer (orange), BASECONV (blue), gradient descent (gray), and least squares (green): we evaluate out-of-distribution ODE parameters (left), forcing function smoothness (middle), and initial condition distribution (right). BASECONV learns a numerically general algorithm that closely matches proper gradient descent and least squares.

1674 D THEORETICAL RESULTS

1675
1676 This section is organized as follows:

- 1677
- 1678 • We detail notation and definitions in Appendix D.1, Appendix D.6.
- 1679
- 1680 • In Appendix D.2, we include theoretical results regarding the primitives from Section 3.3:
- 1681 expressivity results with BASECONV and attention, and iterative algorithms as compositions
- 1682 of primitives.
- 1683 • In Appendix D.3, we discuss upper and lower bounds for implementing gradient descent on
- 1684 least squares using BASECONV, supplementing Section 4.1.
- 1685 • In Appendix D.4, we provide theoretical details regarding BASECONV and polynomials.
- 1686
- 1687 • Finally, in Appendix D.7, we provide technical details about the claims from Section 4.1
- 1688 that BASECONV can perfectly recover Square and Linear.

1689 D.1 NOTATION

1690 We heavily borrow notation from Appendix H of Arora et al. (2023), which we recollect below. We
1691 denote the all 1 row vector of size k , given by $[1 \ 1 \ \dots \ 1 \ 1]$, and the all 0 row vector of size k ,
1692 given by $[0 \ 0 \ \dots \ 0 \ 0]$, as $\mathbf{1}^k$ and $\mathbf{0}^k$, respectively. We also construe the standard basis vector
1693 \mathbf{e}_i as a column vector in this appendix, and adhere to the following matrix indexing convention:
1694 $\mathbf{M}[i, j]$ is the entry in the i th row and the j th column, $\mathbf{M}[i, :] \in \mathbb{F}^{1 \times n}$ denotes the i th row, and
1695 $\mathbf{M}[:, j] \in \mathbb{F}^{m \times 1}$ denotes the j th column of $\mathbf{M} \in \mathbb{F}^{m \times n}$, where \mathbb{F} is a field (the reader can assume
1696 that \mathbb{F} is the field of real numbers i.e. $\mathbb{F} = \mathbb{R}$). We then use $\mathbf{1}^{m \times n}, \mathbf{0}^{m \times n} \in \mathbb{F}^{m \times n}$ to denote the
1697 matrix of all 1s and 0s, respectively. We note that some notation differs from those used in earlier
1698 sections.

1699 Next, we denote the *Hadamard product* of vectors $\mathbf{u}, \mathbf{v} \in \mathbb{F}^n$ as $\mathbf{u} \odot \mathbf{v}$; the operation can be extended
1700 to matrices by applying the Hadamard product column-wise across the matrices. This is commonly
1701 referred to as (*element-wise gating*). For vectors $\mathbf{u}, \mathbf{v} \in \mathbb{F}^n$, we also denote their *linear (or acyclic)*
1702 *convolution* as $\mathbf{u} * \mathbf{v}$ and *cyclic convolution* as $\mathbf{u} \circledast \mathbf{v}$.

1703
1704
1705 **Polynomial Notation.** Since convolution is equivalent to operations on polynomials, it is convenient
1706 to use them to discuss the inputs and outputs of gated convolution models. Let us define maps
1707 $\text{poly} : \mathbb{F}^n \rightarrow \mathbb{F}[X]/(X^n)$ such that

$$1708 \text{poly}(\mathbf{u}) = \sum_{i=0}^{n-1} \mathbf{u}[i]X^i.$$

1709
1710 This allows us to map between vectors and polynomial. Accordingly, we also define $\text{coeff} : \mathbb{F}[X]/(X^{n+1}) \rightarrow \mathbb{F}^n$ as the map converting polynomials back to vectors: $\text{coeff}(\mathbf{u}(X)) = \mathbf{u}$ with
1711 $\mathbf{u}[i]$ defined as the coefficient in $\mathbf{u}(X)$ at degree i .

1712
1713 These operations allow us to interpret the convolution of vectors in terms of polynomial multiplication
1714 (Heideman & Burrus, 1988). More specifically, we have

$$1715 \mathbf{u} * \mathbf{v} = \text{coeff}(\mathbf{u}(X) \cdot \mathbf{v}(X)) \bmod X^n$$

1716
1717 The following notation for a polynomial will be used in this section:

1718
1719 **Definition D.1.** A polynomial $P(X)$ with degree d and some coefficients $\mathbf{c} \in \mathbb{R}^{d+1}$ is defined as,

$$1720 P(X) = \sum_{i=0}^d c_i X^i.$$

1721
1722 Further, the degree of $P(X)$ will be denoted as $\text{deg}(P)$.

Function Approximation. In this part, we collect notation and known results about function approximation. We will reference some definitions from Plešniak (2009); Petersdorff (2015); Smoothness (2006).

The following notation is to denote the k th derivative of a function:

Definition D.2. For some function $f : \mathbb{R} \rightarrow \mathbb{R}$, $f^{(k)} := \frac{d^k}{dx^k} f(x)$ is the k th derivative of f .

Define a set of univariate functions with a notion of continuity:

Definition D.3. We denote $C^k[a, b]$ for $k = 1, 2, \dots$ the space of univariate functions $f : [a, b] \rightarrow \mathbb{R}$, which have derivatives $f^{(1)}, \dots, f^{(k)}$ that are continuous on the closed interval $[a, b]$.

Next we define a set of multivariate functions with a notion of continuity:

Definition D.4. A function $f : [a, b]^n \rightarrow \mathbb{R}$ is in $C^k[a, b]^n$ for $k = 1, 2, \dots$ if all partial derivatives

$$\frac{\partial^\alpha}{\partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \dots \partial x_n^{\alpha_n}} f(y_1, y_2, \dots, y_n)$$

exist and are continuous, for every $\alpha_1, \alpha_2, \dots, \alpha_n \in \mathbb{Z}_{\geq 0}$, such that $\alpha_1 + \alpha_2 + \dots + \alpha_n \leq k$ and every $(y_1, \dots, y_n) \in [a, b]^n$.

We use the following notation for the set of all univariate polynomials:

Definition D.5. For any integer $d \geq 0$, we define

$$\mathcal{P}_d(X) = \{c_0 + c_1 X + \dots + c_d X^d \mid c_k \in \mathbb{R}\}.$$

In other words, $\mathcal{P}_d(X)$ is the space of univariate polynomials of degree less or equal to d .

We use the following notation for multivariate polynomials:

Definition D.6. For any integers $n, d \geq 0$, we define

$$\mathcal{P}_d^n(X_1, \dots, X_n) = \left\{ \sum_{\alpha=(\alpha_1, \dots, \alpha_n) \in \mathbb{Z}_{\geq 0}^n} c_\alpha X_1^{\alpha_1} X_2^{\alpha_2} \dots X_n^{\alpha_n} \mid c_\alpha \in \mathbb{R}, \sum_{i=0}^n \alpha_i \leq d \right\}.$$

Then $\mathcal{P}_d^n(X_1, \dots, X_n)$ is the space of n -variate polynomials of degree less or equal to d .

The following notation is for considering the pointwise absolute value of a matrix:

Definition D.7. For $M \in \mathbb{R}^{N \times D}$ define,

$$\|M\|_\infty = \max_{\substack{0 \leq i < N \\ 0 \leq j < D}} |M[i, j]|.$$

Now lets define the corresponding ∞ -norm for functions:

Definition D.8. For $g : [-1, 1]^{N \times D} \rightarrow \mathbb{R}^{N \times D}$, define

$$\|g\|_\infty = \max_{\mathbf{x} \in [-1, 1]^{N \times D}} |g(\mathbf{x})|.$$

We will use the following version of Jackson's theorem for univariate inputs:

Theorem D.9 (D. Jackson (1930)) Jackson's Theorem for $C^k[-1, 1]$. Let d, k be integers with $d + 1 \geq k \geq 0$ and $f \in C^k[-1, 1]$. Then

$$\inf_{P \in \mathcal{P}_d} \|f - P\|_\infty \leq \left(\frac{\pi}{2}\right)^k \frac{1}{(d+1)d \dots (d-k+2)} \|f^{(k)}\|_\infty. \quad (15)$$

We will use the following version of Jackson's theorem for multivariate inputs:

Theorem D.10 (Plešniak (2009)) Jackson's Theorem for $C^k[-1, 1]^n$. Let d, k be integers with $d + 1 \geq k \geq 0$ and $f \in C^k[-1, 1]^n$. Then

$$\inf_{P \in \mathcal{P}_d^n} \|f - P\|_\infty \leq \frac{c_k}{d^k} \sum_{j=1}^n \left\| \frac{\partial^{k+1}}{\partial x_j^{k+1}} f(\mathbf{x}) \right\|_\infty \quad (16)$$

where c_k is a positive constant.

1782 We will use the following definition of univariate smooth functions:

1783 **Definition D.11.** We call a k times differentiable function $f : [-1, 1] \rightarrow \mathbb{R}$ to be (k, L) -smooth if
 1784 $\|f^{(k)}\|_\infty \leq L$.
 1785

1786 Next, we observe that given a univariate smooth function, there's a univariate bounded degree
 1787 polynomial that approximates it to some error, ϵ :

1788 **Corollary D.12.** For some (k, L) -smooth univariate function f (as in [Definition D.11](#)), then there
 1789 exists a polynomial $P_f(x)$ with
 1790

$$1791 \deg(P_f) \leq O\left(\sqrt[k]{\frac{L}{\epsilon}}\right) + k$$

1792 such that for all $x \in [-1, 1]$

$$1793 |f(x) - P_f(x)| \leq \epsilon.$$

1794 *Proof.* We will be a bit more specific on an upper bound of $\deg(P_f)$. We pick:

$$1795 \deg(P_f) = \left\lceil \frac{\pi}{2} \left(\frac{L}{\epsilon}\right)^{\frac{1}{k}} + k \right\rceil. \quad (17)$$

1800 Let $d = \deg(P_f)$ where P_f is the polynomial that achieves the left hand side of [Equation \(15\)](#). Then
 1801 we have error at most

$$1802 \left(\frac{\pi}{2}\right)^k \frac{1}{(d+1)d \cdots (d-k+2)} \|f^{(k)}\|_\infty.$$

1803 Using the definition of a (k, L) -smooth univariate function in [Definition D.11](#) we get the error at most

$$1804 \left(\frac{\pi}{2}\right)^k \frac{L}{(d+1)d \cdots (d-k+2)} \leq \left(\frac{\pi}{2}\right)^k \frac{L}{(d-k)^k}$$

1805 where the inequality follows since each $d+1, d, \dots, d-k+2 \geq (d-k)$.

1806 Plugging in [Equation \(17\)](#) for d we get the error is at most:

$$1807 \left(\frac{\pi}{2}\right)^k \frac{L}{\left(\frac{\pi}{2}\right)^k \left(\sqrt[k]{\frac{L}{\epsilon}}\right)^k} = \epsilon,$$

1808 as desired. □

1809 We will use the following definition of multivariate smooth functions that map to a single value:

1810 **Definition D.13.** We call a k times differentiable $f : [-1, 1]^n \rightarrow \mathbb{R}$ to be (k, L) -smooth if
 1811 $\left\| \frac{\partial^k}{\partial x_m^k} f(\mathbf{x}) \right\|_\infty \leq L$ for all $1 \leq m \leq n$.
 1812

1813 Now we show the corresponding observation for multivariate functions and polynomials:

1814 **Corollary D.14.** Let $\deg(P_f) = d$. For some (k, L) -smooth multivariate function f (as in [Defini-](#)
 1815 [tion D.13](#)), then there exists a polynomial $P_f(\mathbf{x})$ with
 1816

$$1817 \deg(P_f) \leq O_k\left(\sqrt[k]{\frac{nL}{\epsilon}}\right)$$

1818 such that for all $\mathbf{x} \in [-1, 1]^n$

$$1819 |f(\mathbf{x}) - P_f(\mathbf{x})| \leq \epsilon.$$

1820 *Proof.* Let P_f be the polynomial we get from the left hand side of [Equation \(16\)](#). We want to upper
 1821 bound the error as

$$1822 \frac{c_k}{d^k} \sum_{j=1}^n \left\| \frac{\partial^{k+1}}{\partial x_j^{k+1}} f(\mathbf{x}) \right\|_\infty \leq \epsilon,$$

1836 which follows if

$$1837 \frac{c_k}{d^k} \sum_{j=1}^n L \leq \epsilon$$

1838 since f is (k, L) -smooth. The above is the same as

$$1839 \frac{c_k n L}{d^k} \leq \epsilon,$$

1840 or equivalently

$$1841 \sqrt[k]{\frac{c_k n L}{\epsilon}} \leq d.$$

1842 Picking $d = \left\lceil \sqrt[k]{\frac{c_k n L}{\epsilon}} \right\rceil$ suffices. □

1843 **Arithmetic Circuit Notation.** We briefly recall arithmetic circuits [Peter Bürgisser and Michael Clausen and M. Amin Shokrollah \(1997\)](#). An *arithmetic circuit* \mathcal{C} with variables $X \triangleq \{x_1, x_2, \dots, x_n\}$ over a field \mathbb{F} is interpreted as a directed acyclic graph, where the input nodes are labelled by either the variables from X or constants from \mathbb{F} and the internal nodes are labelled by $+$ or \times with the output being the polynomial computed at the output node.

1844 We shall also refer to the *size*¹ of the circuit \mathcal{C} as the number of wires (or edges in \mathcal{C}), the *depth* of the circuit as the length of the longest path between an input node and the output node, and the *width* of the circuit as the number of wires that will be intersected by a horizontal ‘cut’ through the circuit. Moreover, the *degree* of a circuit is defined as the degree of the polynomial computed by the circuit. We summarize this with the following definition:

1845 **Definition D.15.** An arithmetic circuit \mathcal{C} is an (n, s, Δ, w) -circuit if \mathcal{C} is an n -variate arithmetic circuit of size s , depth at most Δ , and width w .

1846 **BASECONV Architecture.** In the following definitions we formally define the BASECONV model [Arora et al. \(2023\)](#). To formally define BASECONV, we will need the Kaleidoscope hierarchy [Dao et al. \(2020\)](#) as well.

1847 To start, we define butterfly factors:

1848 **Definition D.16.** A **butterfly factor** of size $k \geq 2$ (denoted as $\overline{\mathbf{B}}_k$) is a matrix of the form $\overline{\mathbf{B}}_k = \begin{bmatrix} \mathbf{D}_1 & \mathbf{D}_2 \\ \mathbf{D}_3 & \mathbf{D}_4 \end{bmatrix}$ where each \mathbf{D}_i is a $\frac{k}{2} \times \frac{k}{2}$ diagonal matrix. We restrict k to be a power of 2.

1849 The following definition is for a butterfly factor matrix, which is made up of the above butterfly factors:

1850 **Definition D.17.** A **butterfly factor matrix** of size n with block size k (denoted as $\overline{\mathbf{B}}_k^{(n)}$) is a block diagonal matrix of $\frac{n}{k}$ (possibly different) butterfly factors of size k :

$$1851 \overline{\mathbf{B}}_k^{(n)} = \text{diag} \left([\overline{\mathbf{B}}_k]_1, [\overline{\mathbf{B}}_k]_2, \dots, [\overline{\mathbf{B}}_k]_{\frac{n}{k}} \right)$$

1852 Now lets define a butterfly matrix:

1853 **Definition D.18.** A **butterfly matrix** of size n (denoted as $\overline{\mathbf{B}}^{(n)}$) is a matrix that can be expressed as a product of butterfly factor matrices: $\overline{\mathbf{B}}^{(n)} = \overline{\mathbf{B}}_n^{(n)} \overline{\mathbf{B}}_{\frac{n}{2}}^{(n)} \dots \overline{\mathbf{B}}_2^{(n)}$. Equivalently, we may define $\overline{\mathbf{B}}^{(n)}$ recursively as a matrix that can be expressed in the following form:

$$1854 \overline{\mathbf{B}}^{(n)} = \overline{\mathbf{B}}_n^{(n)} \begin{bmatrix} [\overline{\mathbf{B}}^{(\frac{n}{2})}]_1 & 0 \\ 0 & [\overline{\mathbf{B}}^{(\frac{n}{2})}]_2 \end{bmatrix}$$

1855 (Note that $[\overline{\mathbf{B}}^{(\frac{n}{2})}]_1$ and $[\overline{\mathbf{B}}^{(\frac{n}{2})}]_2$ may be different.)

1856 ¹Note that if all the gates of an arithmetic circuit have bounded arity then the number of wires and gates are asymptotically the same but in this appendix we will consider gates with unbounded arity.

Using these butterfly matrices, let's define the Kaleidoscope Hierarchy:

Definition D.19 (The Kaleidoscope Hierarchy (Dao et al., 2020)).

- Define \mathcal{B} as the set of all matrices that can be expressed in the form $\overline{\mathbf{B}}^{(n)}$ (for some n).
- Define $(\mathcal{B}\mathcal{B}^*)$ as the set of matrices \mathbf{M} of the form $\mathbf{M} = \mathbf{M}_1\mathbf{M}_2^*$ for some $\mathbf{M}_1, \mathbf{M}_2 \in \mathcal{B}$.
- Define $(\mathcal{B}\mathcal{B}^*)^w$ as the set of matrices \mathbf{M} that can be expressed as $\mathbf{M} = \mathbf{M}_w \dots \mathbf{M}_2\mathbf{M}_1$, with each $\mathbf{M}_i \in (\mathcal{B}\mathcal{B}^*)$ ($1 \leq i \leq w$). (The notation w represents width.)
- Define $(\mathcal{B}\mathcal{B}^*)_e^w$ as the set of $n \times n$ matrices \mathbf{M} that can be expressed as $\mathbf{M} = \mathbf{S}\mathbf{E}\mathbf{S}^\top$ for some $en \times en$ matrix $\mathbf{E} \in (\mathcal{B}\mathcal{B}^*)^w$, where $\mathbf{S} \in \mathbb{F}^{n \times en} = [\mathbf{I}_n \ 0 \ \dots \ 0]$ (i.e. \mathbf{M} is the upper-left corner of \mathbf{E}). (The notation e represents expansion relative to n .)

Here we now formally define a BASECONV layer:

Definition D.20 (BASECONV (Arora et al., 2023)). Given an input sequence $\mathbf{u} \in \mathbb{R}^{N \times D}$, where N is the sequence length and D is the model dimension, a learned weight matrix $\mathbf{W} \in \mathbb{R}^{D \times D}$ and biases $\mathbf{B}_1, \mathbf{B}_2 \in \mathbb{R}^{N \times D}$ and a matrix of convolution filters $\mathbf{H} \in \mathbb{R}^{N \times D}$, a BASECONV layer computes the following:

$$\mathbf{y}^{\text{BASECONV}} := (\mathbf{u}\mathbf{W} + \mathbf{B}_1) \odot (\mathbf{H} * \mathbf{u} + \mathbf{B}_2) \in \mathbb{R}^{N \times D}, \quad (18)$$

where the j th column of $\mathbf{H} * \mathbf{u} \in \mathbb{R}^{N \times D}$ is defined as $\mathbf{H}[:, j] * \mathbf{u}[:, j]$.

The corresponding pseudocode for a BASECONV layer is as follows:

Algorithm 1 BASECONV($\mathbf{u}, \mathbf{W}, \mathbf{B}_1, \mathbf{H}, \mathbf{B}_2$)

Require: Input sequence $\mathbf{u} \in \mathbb{R}^{N \times D}$, linear map $\mathbf{W} \in \mathbb{R}^{D \times D}$, convolution filter $\mathbf{H} \in \mathbb{R}^{N \times D}$, and bias matrices $\mathbf{B}_1, \mathbf{B}_2 \in \mathbb{R}^{N \times D}$.

1: In parallel for $0 \leq n < N$: $\mathbf{x}[n, :] = \mathbf{u}[n, :] \cdot \mathbf{W}$

2: In parallel for $0 \leq t < D$: $\mathbf{z}[:, t] = \mathbf{H}[:, t] * \mathbf{u}[:, t]$

3: In parallel for $0 \leq t < D$: $\mathbf{y}[:, t] \leftarrow (\mathbf{x}[:, t] + \mathbf{B}_1[:, t]) \odot (\mathbf{z}[:, t] + \mathbf{B}_2[:, t])$. \triangleright See eq. (18)

4: **return** \mathbf{y}

Remark D.21. The definition of a BASECONV layer in Equation (19) has the input go through a linear layer before the convolution operation. For this section we will assume the linear layer is the identity matrix, as it is not needed for the results in this section.

Assumption D.22. Moving forward we assume the weight matrix $\mathbf{W} \in \mathbb{R}^{D \times D}$ in Definition D.20 also has the property $\mathbf{W} \in (\mathcal{B}\mathcal{B}^*)_{\text{poly-log } D}^{\text{poly-log } D}$. Consequently, each matrix \mathbf{W} has $\tilde{\mathcal{O}}(D)$ parameters and runtime for matrix vector multiplication Dao et al. (2020).

In this section, we will establish some additional basic primitives that we expect need to implement via a BASECONV layer: `shift` and `remember`. We specify them below:

Definition D.23. `shift`(\mathbf{y}, r, t, f)

Shift an sequential input of length N up or down by s entries:

INPUT: $\mathbf{y} \in \mathbb{R}^{N \times D}$, $s \geq 0$.

OUTPUT: $\mathbf{z} \in \mathbb{R}^{N \times D}$ where $\mathbf{z}^+ = \text{shift_down}(\mathbf{y}, s)$ and $\mathbf{z}^- = \text{shift_up}(\mathbf{y}, s)$

$$\begin{array}{c}
1944 \\
1945 \\
1946 \\
1947 \\
1948 \\
1949 \\
1950 \\
1951 \\
1952 \\
1953 \\
1954 \\
1955 \\
1956 \\
1957
\end{array}
\begin{array}{ccc}
\begin{array}{c} \left(\begin{array}{c} \leftarrow \mathbf{y}_0 \rightarrow \\ \vdots \\ \leftarrow \mathbf{y}_{i-1} \rightarrow \\ \leftarrow \mathbf{y}_i \rightarrow \\ \vdots \\ \leftarrow \mathbf{y}_{N-1} \rightarrow \end{array} \right) \\ \mathbf{y} \equiv \end{array} &
\begin{array}{c} \left(\begin{array}{c} \leftarrow \mathbf{0} \rightarrow \\ \vdots \\ \leftarrow \mathbf{0} \rightarrow \\ \leftarrow \mathbf{y}_0 \rightarrow \\ \vdots \\ \leftarrow \mathbf{y}_{N-1-s} \rightarrow \end{array} \right) \\ \mathbf{z}^+ \equiv \end{array} &
\begin{array}{c} \left(\begin{array}{c} \leftarrow \mathbf{y}_s \rightarrow \\ \vdots \\ \leftarrow \mathbf{y}_{N-1} \rightarrow \\ \leftarrow \mathbf{0} \rightarrow \\ \vdots \\ \leftarrow \mathbf{0} \rightarrow \end{array} \right) \\ \mathbf{z}^- \equiv \end{array}
\end{array}$$

The following proposition is defining the convolution Kernel that computes the $\text{shift_down}(\cdot, \lfloor \frac{N}{2} \rfloor)$ primitive:

Proposition D.24. Define $\mathbf{H} \in \mathbb{R}^{2N \times D}$ as

$$\mathbf{H}[k, :] = \begin{cases} \mathbf{1}^D & \text{if } k = N \\ 0 & \text{otherwise} \end{cases}.$$

For any $\mathbf{u} \in \mathbb{R}^{2N \times D}$, $\mathbf{H} * \mathbf{u}$ will result in

$$\mathbf{H} * \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix} \rightarrow \begin{pmatrix} \mathbf{0}^{N \times D} \\ \mathbf{u}_1 \end{pmatrix},$$

where $\mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^{N \times D}$.

Proof. The convolution operation: $\mathbf{H} * \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix}$ where each column of \mathbf{H} is convolved with each column of \mathbf{u} can be restated as a polynomial multiplication. For column i , $0 \leq i < 2N$,

$$\mathbf{H}[:, i] * \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix}[:, i] = \text{coeff}((X^N \cdot \mathbf{u}[:, i](X)) \bmod X^{2N}).$$

Note that the columns of \mathbf{H} are all \mathbf{e}_N basis vectors and $\text{poly}(\mathbf{e}_N) = X^N$.

When we multiply the term through the input polynomial we get,

$$\begin{aligned}
& \text{coeff}(X^N \cdot (\mathbf{u}[0][i] + \mathbf{u}[1][i]X + \dots + \mathbf{u}[2N-1][i]X^{2N-1}) \bmod X^{2N}) \\
& = \text{coeff}(\mathbf{u}[0][i]X^N + \mathbf{u}[1][i]X^{N+1} + \dots + \mathbf{u}[2N-1][i]X^{3N-1} \bmod X^{2N}).
\end{aligned}$$

With the lower order terms all becoming zeros, the above is same as

$$\begin{aligned}
& \text{coeff}((0 + 0X + \dots + 0X^{N-1} \\
& + \mathbf{u}[0][i]X^N + \mathbf{u}[1][i]X^{N+1} + \dots + \mathbf{u}[2N-1][i]X^{3N-1}) \bmod X^{2N}).
\end{aligned}$$

After we take the $\bmod X^{2N}$ we get

$$\text{coeff}(0 + 0X + \dots + 0X^{N-1} + \mathbf{u}[0][i]X^N + \dots + \mathbf{u}[N-1][i]X^{2N-1}),$$

which implies that $\mathbf{H} * \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{pmatrix}$ is

$$\begin{pmatrix} \mathbf{0}^{N \times D} \\ \mathbf{u}_1 \end{pmatrix},$$

as desired. \square

We also define the following primitive:

Definition D.25. $\text{remember}(\mathbf{y}, r, t, f)$

INPUT: $\mathbf{y} \in \mathbb{R}^{N' \times d'}$, $r \in \mathbb{Z}, t \in \mathbb{Z}, f: \mathbb{R}^{t-r} \rightarrow \mathbb{R}^{t-r+s}$, $\mathbf{v}_1 \in \mathbb{R}^r, \mathbf{x} \in \mathbb{R}^{t-r}$, where \mathbf{y} is defined as below.

OUTPUT: $\mathbf{z} \in \mathbb{R}^{N' \times d'}$, which is defined as follows:

$$\mathbf{y} \equiv \begin{pmatrix} \overleftarrow{\mathbf{v}_1} \overrightarrow{\phantom{\mathbf{v}_1}} \\ \overleftarrow{\mathbf{x}} \overrightarrow{\phantom{\mathbf{x}}} \\ \mathbf{0}^{s \times d'} \\ \overleftarrow{\mathbf{v}_2} \overrightarrow{\phantom{\mathbf{v}_2}} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix} \quad \mathbf{z} \equiv \begin{pmatrix} \overleftarrow{\mathbf{v}_1} \overrightarrow{\phantom{\mathbf{v}_1}} \\ \overleftarrow{f(\mathbf{x})} \overrightarrow{\phantom{f(\mathbf{x})}} \\ \overleftarrow{\mathbf{v}_2} \overrightarrow{\phantom{\mathbf{v}_2}} \\ \mathbf{0} \\ \vdots \\ \mathbf{0} \end{pmatrix}$$

We will need the following BASECONV implementation of `remember`:

Proposition D.26 (Arora et al. (2024), The Remembering Primitive). *For any $\mathbf{x} \in \mathbb{R}^{n \times d'}$, $\mathbf{v}_1 \in \mathbb{R}^{r \times d'}$, $\mathbf{v}_2 \in \mathbb{R}^{m-r}$ where $n = t - r$ contained in some $\mathbf{y} \in \mathbb{R}^{N' \times d'}$ such that \mathbf{v}_1 is in the first r rows, \mathbf{x} is in the next n rows, $\mathbf{0}$ s fill up the next s rows, and \mathbf{v}_2 are in the next $m - r$ rows, for some $3n + 3m + 2s + 2t \leq N'$ so that for $\mathbf{h} \in \mathbb{R}^{n \times d}$ and $\mathbf{W} \in \mathbb{R}^{d' \times d}$ with $\mathbf{x} * \mathbf{h} \in \mathbb{R}^{(n+s) \times d'}$ and $\mathbf{v} * \mathbf{h} \in \mathbb{R}^{(m+t) \times d'}$, where $\mathbf{v} \in \mathbb{R}^{m \times d'}$ is defined as $\mathbf{v}_2 + \text{shift_down}(\mathbf{v}_1, m - r)$, there exists a $(N', 8, d', N', d')$ – BASECONV that computes $\text{remember}(\mathbf{y}, r, t, f)$, where f can be implemented in 1 layer of BASECONV through the parameters $\mathbf{W} \in \mathbb{R}^{d' \times d}$, $\mathbf{h} \in \mathbb{R}^{n \times d}$, $\mathbf{b}_1 \in \mathbb{R}^{s \times d'}$, $\mathbf{b}_2 \in \mathbb{R}^{(m+t) \times d'}$ as defined below:*

$$f(\mathbf{u}) = \left(\begin{pmatrix} \mathbf{u}\mathbf{W} \\ \mathbf{0}^{s \times d'} \end{pmatrix} + \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{1}^{s \times d'} \end{pmatrix} \right) \odot \left(\mathbf{u} * \mathbf{h} + \begin{pmatrix} \mathbf{b}_2 \\ \mathbf{0}^{s \times d'} \end{pmatrix} \right)$$

We will also need the following generalization of the above result:

Corollary D.27 (Arora et al. (2023)). *Let \mathbf{y} be as in Proposition D.26 but now let f be implemented with BASECONV(N, L, D, N, D). Then $\text{remember}(\mathbf{y}, r, t, f)$ where $t - r = n$ can be implemented with BASECONV via $(N, O(L), D, N, D)$ – BASECONV.*

The rest of Appendix D will use this 5–tuple notation for BASECONV:

Definition D.28. Lets define a 5-tuple notation for a BASECONV layer as (N, ℓ, D, N', D') – BASECONV with ℓ layers such that:

1. Input and output are $N \times D$ matrices.
2. Each layer is defined by Definition D.20 where N and D are replaced by N' and D' . I.e. each layer takes in $N' \times D'$ matrices and output $N' \times D'$ matrices. We refer to the tuple (N', D') as the *inner dimension* of the model.
3. The matrices are projected from $(N, D) \rightarrow (N', D')$ (and vice-versa) via a linear projection.

We state the following bounds on parameters and runtime for a single BASECONV layer:

Proposition D.29 (Arora et al. (2023)). *An $(N, 1, D, N, D)$ – BASECONV requires $\tilde{O}(ND)$ parameters and runtime.*

We state the following result that says arithmetic circuit can be represented as a BASECONV model:

Theorem D.30 (Arora et al. (2023), Theorem H.21). *For any (ND, s, Δ, w) -arithmetic circuit \mathcal{C} , there exists an equivalent (N, Δ', D, N', D') – BASECONV with $\Delta' = \mathcal{O}(\Delta \log w)$, $N' = \mathcal{O}(w)$, $D' = D$ that simulates \mathcal{C} .*

D.2 PRIMITIVES

In this section, we provide theoretical results about primitives.

- In Appendix D.2.1, we implement the three primitives (READ, LINEAR, and MULTIPLY) from Section 3.3 using BASECONV, each using a single layer.
- Next, in Appendix D.2.2 and D.2.3, we briefly sketch how the three primitives READ, LINEAR, and MULTIPLY can be used in composition to exactly express gradient descent and Newton’s method iterations on least squares (see Appendix A).
- Finally, in Appendix D.2.4, we provide a proof that a single layer of causal softmax attention cannot exactly represent the entry-wise squaring function. As a corollary, since entry-wise square is a special case of MULTIPLY, this implies that attention cannot exactly express the MULTIPLY task for all arguments.

BASECONV parameterization We recount the parameterization of BASECONV from Equation 2:

$$\begin{aligned} \mathbf{y} &:= \left(\underbrace{(\mathbf{u} \cdot \mathbf{W}_{gate} + \mathbf{b}_{gate})}_{\text{Linear Projection}} \odot \underbrace{(\mathbf{h} * (\mathbf{u} \cdot \mathbf{W}_{in} + \mathbf{b}_{in}) + \mathbf{b}_{conv})}_{\text{Convolution}} \right) \cdot \mathbf{W}_{out} + \mathbf{b}_{out} \\ &:= W_{out}(W_{gate}(\mathbf{u}) \odot Conv(W_{in}(\mathbf{u}))) \end{aligned} \quad (19)$$

where $W_{in}, W_{gate}, W_{out}$ are linear projections $\mathbb{R}^D \rightarrow \mathbb{R}^D$.

D.2.1 1-LAYER BASECONV CAN IMPLEMENT LINEAR ALGEBRA PRIMITIVES

Below, we formally define the linear algebra primitives we discuss in Section 3.3, and we describe our BASECONV weight constructions.

Read The READ operator, which maps inputs $\mathbf{u} \in \mathbb{R}^{N \times d}$ to outputs $\mathbf{y} \in \mathbb{R}^{N \times d}$, is:

$$\text{READ}(i, j, a, b)(\mathbf{u}) = \begin{cases} \mathbf{u}[k, a : b] & k \neq j \\ \mathbf{u}[i, a : b] & k = j \end{cases}. \quad (20)$$

Our implementation requires the use of the positional encodings and residual connections within the BASECONV architecture. Concretely, consider the input

$$\mathbf{u}_{in} = \begin{pmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \dots & \mathbf{e}_N \\ \mathbf{u}[1, :] & \mathbf{u}[2, :] & \dots & \mathbf{u}[N, :] \end{pmatrix},$$

where the basis vector e_k represents the positional encoding for the k -th entry of the sequence. Define the output of the BASECONV layer *with residual connection*:

$$\mathbf{y} := W_{out}(W_{gate}(\mathbf{u}) \odot Conv(W_{in}(\mathbf{u})) + \mathbf{u}).$$

Then the following weight construction is equivalent to $\text{READ}(i, j, a, b)$:

- $W_{gate}(\mathbf{u}[k, :]) := \mathbf{u}[k, j] \mathbf{1}^D$
- $Conv(W_{in}(\mathbf{u}))[k, :] := \mathbf{u}[k + i - j, :] - \mathbf{u}[k, :]$
- $W_{out} := \text{proj}(a : b)$.

In particular, W_{gate} is defined such that

$$W_{gate}(\mathbf{u}[k, :]) = \begin{cases} \mathbf{1}^D & k = j \\ \mathbf{0}^D & k \neq j \end{cases}.$$

Thus

$$W_{gate}(\mathbf{u}) \odot Conv(W_{in}(\mathbf{u})) = \begin{cases} \mathbf{u}[k + i - j, :] - \mathbf{u}[k, :] = \mathbf{u}[i, :] - \mathbf{u}[j, :] & k = j \\ \mathbf{0}^D & k \neq j \end{cases}.$$

Finally,

$$W_{gate}(\mathbf{u}) \odot Conv(W_{in}(\mathbf{u})) + \mathbf{u} = \begin{cases} \mathbf{u}[i, :] & k = j \\ \mathbf{u}[k, :] & k \neq j \end{cases}$$

so the final output of this layer will be exactly equivalent to $\text{READ}(i, j, a, b)$.

2106 **Linear transformation** The LINEAR operator, which maps inputs $\mathbf{u} \in \mathbb{R}^{N \times d_{in}}$ to outputs $\mathbf{y} \in$
 2107 $\mathbb{R}^{N \times d_{out}}$, is:

$$2108 \text{LINEAR}(\mathbf{H})(\mathbf{u}) = \mathbf{u}\mathbf{H} \quad (21)$$

2109 where $\mathbf{H} : \mathbb{R}^{d_{in}} \rightarrow \mathbb{R}^{d_{out}}$ is a linear map.

2110 Define $\text{Conv}(W_{in}(\mathbf{u})) = \mathbf{1}_D$, $W_{gate} = I$, and $W_{out} = \mathbf{H}$. Then

$$2111 W_{gate}(\mathbf{u}) \odot \text{Conv}(W_{in}(\mathbf{u})) = \mathbf{u}$$

2112 so

$$2113 W_{out}(W_{gate}(\mathbf{u}) \odot \text{Conv}(W_{in}(\mathbf{u}))) = \mathbf{u}\mathbf{H}.$$

2114 Thus the output of this layer is exactly equivalent to $\text{LINEAR}(\mathbf{H})$.

2115 **Element-wise multiply** The MULTIPLY operator, which maps inputs $\mathbf{u} \in \mathbb{R}^{N \times d_{in}}$ to outputs
 2116 $\mathbf{y} \in \mathbb{R}^{N \times d_{out}}$, is:

$$2117 \text{MULTIPLY}(a, b, d_{out})(\mathbf{u}) = \mathbf{u}[:, a : a + d_{out}] \odot \mathbf{u}[:, b : b + d_{out}] \quad (22)$$

2118 Define $\text{Conv} = \text{Identity}$, $W_{in} = \text{proj}(a : a + d_{out})$, $W_{gate} = \text{proj}(b : b + d_{out})$, and $W_{out} = I$.

2119 Then

$$2120 W_{gate}(\mathbf{u}) \odot \text{Conv}(W_{in}(\mathbf{u})) = \mathbf{u}[:, a : a + d_{out}] \odot \mathbf{u}[:, b : b + d_{out}].$$

2121 Since $W_{out} = I$, the output of this layer will be equivalent to $\text{MULTIPLY}(a, b, d_{out})$.

2122 D.2.2 GRADIENT DESCENT

2123 We assume our input is of the form

$$2124 \mathbf{u} = \begin{pmatrix} \mathbf{a}_1 & \dots & \mathbf{a}_N & \mathbf{x}_0 \\ b_1 & \dots & b_N & 0 \end{pmatrix}.$$

2125 Our goal is to compute the gradient update

$$2126 \mathbf{x}_1 := \mathbf{x}_0 - \eta \sum_{i=1}^N (\mathbf{x}_0^T \mathbf{a}_i - b_i) \mathbf{a}_i. \quad (23)$$

2127 Intuitively, our argument proceeds similarly to the causal gradient descent construction from Ap-
 2128 pendix D.3.1:

- 2129 • First, we repeatedly apply READ and LINEAR to move the information $\{\mathbf{a}_i, b_i\} \forall i$ into e.g.
 2130 the final entry of the sequence. Without loss of generality, we omit the rest of the sequence,
 2131 and assume we have access to a large enough embedding dimension that we can make use
 2132 of arbitrary amounts of memory.

2133 After this phase, our \mathbf{u} is of the form

$$2134 \dots (\mathbf{x}_0 \ 0 \ \mathbf{a}_1 \ \dots \ \mathbf{a}_N \ b_1 \ \dots \ b_N \ \dots)^T.$$

- 2135 • Next, we use MULTIPLY and LINEAR to compute and store $\{\mathbf{x}_0^T \mathbf{a}_i\}$ for all i . We will end
 2136 up with

$$2137 \mathbf{u} = \dots (\mathbf{x}_0 \ 0 \ \{\mathbf{a}_i\}_i \ \{b_i\}_i \ \{\mathbf{x}_0^T \mathbf{a}_i\}_i \ \dots).$$

- 2138 • We use LINEAR to compute and store $\{\mathbf{x}_0^T \mathbf{a}_i - b_i\}$ for all i :

$$2139 \mathbf{u} = \dots (\mathbf{x}_0 \ 0 \ \{\mathbf{a}_i\}_i \ \{b_i\}_i \ \{\mathbf{x}_0^T \mathbf{a}_i\}_i \ \{\mathbf{x}_0^T \mathbf{a}_i - b_i\}_i \ \dots).$$

- 2140 • We use MULTIPLY and LINEAR to compute and store $\{(\mathbf{x}_0^T \mathbf{a}_i - b_i) \mathbf{a}_i\}$ for all i :

$$2141 \mathbf{u} = \dots (\mathbf{x}_0 \ 0 \ \{\mathbf{a}_i\}_i \ \{b_i\}_i \ \{\mathbf{x}_0^T \mathbf{a}_i\}_i \ \{(\mathbf{x}_0^T \mathbf{a}_i - b_i) \mathbf{a}_i\}_i \ \dots).$$

- 2142 • Finally, we can use LINEAR to compute the gradient update:

$$2143 \mathbf{u} = \dots (\mathbf{x}_0 - \eta \sum_{i=1}^N (\mathbf{x}_0^T \mathbf{a}_i - b_i) \mathbf{a}_i \ 0 \ \{\mathbf{a}_i\}_i \ \{b_i\}_i \ \{\mathbf{x}_0^T \mathbf{a}_i\}_i \ \{(\mathbf{x}_0^T \mathbf{a}_i - b_i) \mathbf{a}_i\}_i \ \dots).$$

2160 D.2.3 NEWTON’S METHOD
2161

2162 We assume our input is of the form

$$2163 \mathbf{u} = \begin{pmatrix} \mathbf{a}_1 & \dots & \mathbf{a}_N & \mathbf{M}_0[1, :] & \dots & \mathbf{M}_0[D, :] \\ b_1 & \dots & b_N & 0 & \dots & 0 \end{pmatrix}. \quad 2164$$

2165 Our goal is to compute the Newton’s iterate:

$$2166 \mathbf{M}_1 := \mathbf{M}_0(2\mathbf{I} - (\mathbf{a}^T \mathbf{a})\mathbf{M}_0), \quad 2167 \quad (24)$$

2168 where

$$2169 \mathbf{a} = \begin{pmatrix} \leftarrow \mathbf{a}_1 \rightarrow \\ \vdots \\ \leftarrow \mathbf{a}_N \rightarrow \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} b_1 \\ \vdots \\ b_N \end{pmatrix}. \quad 2170 \quad (25)$$

2171 For any matrix $\mathbf{M} \in \mathbb{R}^{n \times p}$, let *flt* denote the `flatten` operation, so that *flt*(\mathbf{M}) represent a
2172 vectorized version of \mathbf{M} : $\text{flt}(\mathbf{M}) \in \mathbb{R}^{np}$.

2173 We proceed similarly to the argument from Appendix D.2.2.

- 2174 • First, we repeatedly apply `READ` and `LINEAR` to move all information $\{\mathbf{a}_i\}_i \forall i$ and $\text{flt}(\mathbf{M})$
2175 to e.g. the final entry of the sequence. We omit the rest of the sequence for notational ease,
2176 and we assume we have access to a large enough embedding dimension that we can make
2177 use of arbitrary amounts of memory.

2178 After this phase, we have

$$2179 \mathbf{u} = \dots (\text{flt}(\mathbf{M}_0) \quad \{\mathbf{a}_i\}_i \quad \dots). \quad 2180$$

- 2181 • Using `LINEAR`, we can copy and rearrange the \mathbf{a}_i ’s to construct copies of $\text{flt}(\mathbf{a})$ and
2182 $\text{flt}(\mathbf{a}^T)$:

$$2183 \mathbf{u} = \dots (\text{flt}(\mathbf{M}_0) \quad \{\mathbf{a}_i\}_i \quad \text{flt}(\mathbf{a}^T) \quad \text{flt}(\mathbf{a}) \quad \dots). \quad 2184$$

- 2185 • Now, note that we can represent the matrix multiplication $\mathbf{a}^T \mathbf{a}$ as a linear combination of
2186 the entries of the element-wise multiplication $\text{flt}(\mathbf{a}^T) \odot \text{flt}(\mathbf{a})$. This means that we can
2187 obtain $\text{flt}(\mathbf{a}^T \mathbf{a})$ using a single application of `MULTIPLY` and `LINEAR`:

$$2188 \mathbf{u} = \dots (\text{flt}(\mathbf{M}_0) \quad \{\mathbf{a}_i\}_i \quad \text{flt}(\mathbf{a}^T) \quad \text{flt}(\mathbf{a}) \quad \text{flt}(\mathbf{a}^T \mathbf{a}) \dots). \quad 2189$$

- 2190 • By the same argument, we can obtain $\text{flt}((\mathbf{a}^T \mathbf{a})\mathbf{M}_0)$ using another application of `MULTI-`
2191 `PPLY` and `LINEAR`:

$$2192 \mathbf{u} = \dots (\text{flt}(\mathbf{M}_0) \quad \{\mathbf{a}_i\}_i \quad \text{flt}(\mathbf{a}^T) \quad \text{flt}(\mathbf{a}) \quad \text{flt}((\mathbf{a}^T \mathbf{a})\mathbf{M}_0) \dots). \quad 2193$$

- 2194 • Finally, we have that $\text{flt}(\mathbf{M}_1) := 2\text{flt}(\mathbf{M}_0) - \text{flt}((\mathbf{a}^T \mathbf{a})\mathbf{M}_0)$ can be obtained using
2195 `LINEAR` once more:

$$2196 \mathbf{u} = \dots (\text{flt}(\mathbf{M}_1) \quad \{\mathbf{a}_i\}_i \quad \text{flt}(\mathbf{a}^T) \quad \text{flt}(\mathbf{a}) \quad \text{flt}((\mathbf{a}^T \mathbf{a})\mathbf{M}_0) \dots). \quad 2197$$

2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

2214 D.2.4 SOFTMAX ATTENTION CAN'T IMPLEMENT ELEMENT-WISE SQUARING.
 2215

2216 In this section, we consider the following parameterization of *softmax attention*:

$$2217 \text{Attn}(\mathbf{u}) = \text{softmax}((\mathbf{u}\mathbf{W}_Q)(\mathbf{u}\mathbf{W}_K)^T + \mathbf{M})(\mathbf{u}\mathbf{W}_V + \mathbf{B}), \quad (26)$$

2219 where $\mathbf{u} \in \mathbb{R}^{N \times D}$, $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{D \times D}$, $\mathbf{B} \in \mathbb{R}^{N \times D}$, and $\mathbf{M} \in \mathbb{R}^{N \times N}$ is the causal attention
 2220 mask:

$$2221 \mathbf{M}_{ij} = \begin{cases} -\infty & i < j \\ 0 & \text{otherwise} \end{cases} \quad (27)$$

2224 **Theorem D.31.** *One-layer single-headed causal softmax attention cannot exactly represent the*
 2225 *entry-wise squaring function* $\text{SQUARE} : \mathbb{R}^{N \times D} \rightarrow \mathbb{R}^{N \times D}$ *s.t.*

$$2226 \text{SQUARE}(\mathbf{u})_{ij} = \mathbf{u}_{ij}^2$$

2227 *for all* $\mathbf{u} \in \mathbb{R}^{N \times D}$.
 2228

2231 *Proof.* We proceed by contradiction. Let's assume there exists $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V, \mathbf{B} \in \mathbb{R}^{D \times D}$ and
 2232 $\mathbf{M} \in \mathbb{R}^{N \times N}$ such that $\forall \mathbf{u} \in \mathbb{R}^{N \times D}$,

$$2233 \text{softmax}((\mathbf{u}\mathbf{W}_Q)(\mathbf{u}\mathbf{W}_K)^T + \mathbf{M})(\mathbf{u}\mathbf{W}_V + \mathbf{B}) = \text{SQUARE}(\mathbf{u}). \quad (28)$$

2234 Consider the set of inputs $\mathbf{u} \in \mathbb{R}^{N \times D}$ with at most one non-zero entry, defined as

$$2235 \mathbf{u}_{ij} = \begin{cases} \mathbf{u}_{ij} & (i, j) = (a, b) \\ 0 & \text{else} \end{cases} \quad (29)$$

2236 for an arbitrary choice of $a \in [N], b \in [D]$. Then:

$$2237 \mathbf{Q} := \mathbf{u}\mathbf{W}_Q = \begin{pmatrix} \mathbf{0}^N \\ \vdots \\ \mathbf{0}^N \\ \mathbf{u}_{ab}\mathbf{W}_Q[b, :] \\ \vdots \\ \mathbf{0}^N \end{pmatrix} \quad (30)$$

2238 where \mathbf{Q} 's rows are all $\mathbf{0}$ except for the a -th, which is $\mathbf{u}_{ab}\mathbf{W}_Q[b, :]$.

2239 Similarly:

$$2240 \mathbf{K} := \mathbf{u}\mathbf{W}_K = \begin{pmatrix} \mathbf{0}^N \\ \vdots \\ \mathbf{0}^N \\ \mathbf{u}_{ab}\mathbf{W}_K[b, :] \\ \vdots \\ \mathbf{0}^N \end{pmatrix} \quad (31)$$

2268 and

$$2269 \quad \mathbf{V} := \mathbf{u}\mathbf{W}_V = \begin{pmatrix} \mathbf{0}^N \\ 2270 \quad \vdots \\ 2271 \quad \mathbf{0}^N \\ 2272 \quad \mathbf{u}_{ab}\mathbf{W}_V[b, :] \\ 2273 \quad \vdots \\ 2274 \quad \mathbf{0}^N \end{pmatrix} \quad (32)$$

2282 Then the pre-softmax attention matrix, $\mathbf{A}' = \mathbf{Q}\mathbf{K}^T$, satisfies

$$2283 \quad \mathbf{A}'_{ij} = \begin{cases} \mathbf{u}_{ab}^2(\mathbf{W}_Q\mathbf{W}_K^T)_{bb} & (i, j) = (a, a) \\ 2284 \quad 0 & \text{otherwise} \end{cases}. \quad (33)$$

2286 Define

$$2287 \quad C := \mathbf{u}_{ab}^2(\mathbf{W}_Q\mathbf{W}_K^T)_{bb}. \quad (34)$$

2289 Now consider what happens after we apply the softmax operator. Recall that the softmax operator is defined as

$$2291 \quad \text{softmax}(\mathbf{z})_i = \frac{\exp(\mathbf{z}_i)}{\sum_{j=1}^D \exp(\mathbf{z}_j)} \quad (35)$$

2293 for $\mathbf{z} \in \mathbb{R}^D$. Then $\mathbf{A} := \text{softmax}(\mathbf{A}' + \mathbf{M})$ satisfies

$$2295 \quad \mathbf{A}_{ij} = \begin{cases} \frac{1}{i} & i \neq a \\ 2296 \quad \frac{1}{\exp(C)+a-1} & i = a, j \neq a \\ 2297 \quad \frac{\exp(C)}{\exp(C)+a-1} & (i, j) = (a, a) \end{cases} \quad (36)$$

2300 Now let's consider the output of softmax attention:

$$2301 \quad \mathbf{O} = \mathbf{A}(\mathbf{V} + \mathbf{B}) \quad (37)$$

2303 such that $\mathbf{O} = \text{SQUARE}(\mathbf{u})$.

2304 Note that for $i \neq a$:

$$2306 \quad \mathbf{O}[i, :] = \frac{1}{i} \sum_{k=1}^i (\mathbf{V} + \mathbf{B})[k, :] \quad (38)$$

2308 and this must also be equal to $\mathbf{0}^N = \text{SQUARE}(\mathbf{u})[i, :]$. We consider three cases:

- 2311 • First, consider $i < a$ in order from $i = 1, \dots, a - 1$. Since this equality is true for all $i < a$,
- 2312 we can verify that $(\mathbf{V} + \mathbf{B})[i, :]$ must equal $\mathbf{0}^N$ for all $i < a$.

- 2313 • Next, looking at $i = a + 1$, we have

$$2315 \quad \frac{1}{a+1} ((\mathbf{V} + \mathbf{B})[a, :] + (\mathbf{V} + \mathbf{B})[a+1, :]) = \mathbf{0}^N \quad (39)$$

2317 so we must have

$$2318 \quad (\mathbf{V} + \mathbf{B})[a, :] = -(\mathbf{V} + \mathbf{B})[a+1, :] \quad (40)$$

- 2319 • Finally, from $i \geq a + 1$, we can again conclude that $(\mathbf{V} + \mathbf{B})[i, :]$ must equal $\mathbf{0}^N$ for all
- 2321 $i > a + 1$.

This means the only rows of $\mathbf{V} + \mathbf{B}$ that might not be zero are $(\mathbf{V} + \mathbf{B})[a, :]$ and $(\mathbf{V} + \mathbf{B})[a + 1, :]$. Thus looking at the a -th row:

$$\begin{aligned} \frac{\exp(C)}{\exp(C) + a - 1} (\mathbf{V} + \mathbf{B})[a, :] &= \text{SQUARE}(\mathbf{u})[a, :] \\ &= [0 \quad \dots \quad \mathbf{u}_{ab}^2 \quad \dots \quad 0] \end{aligned}$$

Recall that from above,

$$\mathbf{V}[a, :] = \mathbf{u}_{ab} \mathbf{W}_V[b, :] \quad (41)$$

Then analyzing entry-wise, we have:

$$\frac{\exp(C)}{\exp(C) + a - 1} (\mathbf{u}_{ab} \mathbf{W}_V[b, j] + \mathbf{B}[a, j]) = 0 \quad (42)$$

for all $j \neq b$, and

$$\frac{\exp(C)}{\exp(C) + a - 1} (\mathbf{u}_{ab} \mathbf{W}_V[b, b] + \mathbf{B}[a, b]) = \mathbf{u}_{ab}^2. \quad (43)$$

We now plug back in our expression for C and simplifying the latter equation. For ease of notation, denote $A := (\mathbf{W}_Q \mathbf{W}_K^T)_{bb}$, $V := \mathbf{W}_V[b, b]$, and $B := \mathbf{B}[a, b]$. Then the expression simplifies to:

$$V \exp(A \mathbf{u}_{ab}^2) \mathbf{u}_{ab} + B \exp(A \mathbf{u}_{ab}^2) = \exp(A \mathbf{u}_{ab}^2) \mathbf{u}_{ab}^2 + (a - 1) \mathbf{u}_{ab}^2$$

This must hold for all non-zero values of \mathbf{u}_{ab} . We can take $V = B = 0$, but we are still left with

$$\begin{aligned} - \exp(A \mathbf{u}_{ab}^2) \mathbf{u}_{ab}^2 &= (a - 1) \mathbf{u}_{ab}^2 \\ - \exp(A \mathbf{u}_{ab}^2) &= a - 1 \end{aligned}$$

However, there is no choice of A such that this statement holds. This completes the proof by contradiction.

As a corollary, we have

Corollary D.32. *One-layer single-headed causal softmax attention cannot exactly represent the entry-wise multiply function $\text{MULTIPLY} : \mathbb{R}^{N \times D} \rightarrow \mathbb{R}^{N \times d_{out}}$ s.t.*

$$\text{MULTIPLY}(a, b, d_{out})(\mathbf{u}) = \mathbf{u}[:, a : a + d_{out}] \odot \mathbf{u}[:, b : b + d_{out}] \quad (44)$$

for all $\mathbf{u} \in \mathbb{R}^{N \times D}$ and all choices of a, b, d_{out} .

Proof. Note that for $a = 0, b = 0$, and $d_{out} = D$,

$$\text{SQUARE}(\mathbf{u}) = \text{MULTIPLY}(a, b, d_{out})(\mathbf{u}).$$

Since softmax attention cannot exactly represent SQUARE for all \mathbf{u} , it also cannot represent MULTIPLY for all \mathbf{u} . \square

2376 D.3 UPPER AND LOWER BOUNDS WITH BASECONV FOR GRADIENT DESCENT
2377

2378 In this section, we detail upper and lower bounds for implementing gradient descent using BASECONV, as discussed in Section 4.1.
2379

- 2380
- 2381 • **Upper bounds.** We provide two explicit constructions for implementing iterations gradient
2382 descent on linear regression: one for *non-causal* BASECONV requiring $O(1)$ layers and
2383 $O(D)$ state size, and one for *causal* BASECONV requiring $O(1)$ layers and $O(D^2)$ state
2384 size.
 - 2385 • **Lower bounds.** In Appendix D.3.2, we prove that our constructions are asymptotically
2386 optimal with respect to layers and state size.
2387

2388 D.3.1 UPPER BOUNDS: BASECONV CAN IMPLEMENT GRADIENT DESCENT FOR LINEAR
2389 REGRESSION
2390

2391 In this section, we provide weight constructions for exactly implementing gradient descent on linear
2392 regression. Recall:

2393
$$\mathcal{L}_N = \frac{1}{2N} \sum_{i=1}^N (\mathbf{x}^T \mathbf{a}_i - \mathbf{b}_i)^2 \quad (45)$$

2394
2395

2396 so

2397
$$\nabla_{\mathbf{x}} \mathcal{L}_N = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^T \mathbf{a}_i - \mathbf{b}_i) \mathbf{a}_i \quad (46)$$

2398
2399

2400
$$= \frac{1}{N} \left(\sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i - \left(\sum_{i=1}^N \mathbf{a}_i \mathbf{a}_i^T \right) \mathbf{x} \right) \quad (47)$$

2401
2402
2403

2404 **Non-causal BASECONV** This weight construction uses Equation 46 to compute the gradient
2405 descent update.

2406 We note that non-causal constructions for in-context linear regression are standard in the literature:
2407 e.g. Von Oswald et al. (2023); Ahn et al. (2024).
2408

2409 We start with input:

2410
$$\mathbf{b} \equiv \begin{pmatrix} \mathbf{a}_1 & \dots & \mathbf{a}_N & \mathbf{a}_q \\ \mathbf{b}_1 & \dots & \mathbf{b}_N & 0 \end{pmatrix}$$

2411
2412

2413 We define the initial embedding:
2414

2415
$$\begin{pmatrix} \mathbf{a}_1 & \dots & \mathbf{a}_N & \mathbf{0}^D \\ \mathbf{b}_1 & \dots & \mathbf{b}_N & 0 \\ \mathbf{x}_0 & \dots & \mathbf{x}_0 & \mathbf{x}_0 \\ \mathbf{0}^D & \dots & \mathbf{0}^D & \mathbf{0}^D \\ \mathbf{0}^D & \dots & \mathbf{0}^D & \mathbf{0}^D \\ \mathbf{0}^D & \dots & \mathbf{0}^D & \mathbf{a}_q \\ 0 & \dots & 0 & 0 \end{pmatrix}$$

2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428

2429 We drop the bottom two rows of the block matrix representation for now and show how to perform the gradient descent update with the rest of the embedding.

2430 Layer 1:

$$\begin{array}{l}
 2431 \\
 2432 \\
 2433 \\
 2434 \\
 2435 \\
 2436 \\
 2437 \\
 2438 \\
 2439 \\
 2440 \\
 2441 \\
 2442 \\
 2443 \\
 2444 \\
 2445 \\
 2446 \\
 2447 \\
 2448 \\
 2449 \\
 2450 \\
 2451
 \end{array}$$

$$\begin{array}{c}
 \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{0}^D \rightarrow \end{array} \right) \odot \left(\begin{array}{c} \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1} \rightarrow \\ \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{0}^D \rightarrow \end{array} \right) = \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \odot \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{0}^D \rightarrow \end{array} \right) \\
 \underbrace{\hspace{10em}}_{\text{conv}(\text{in_proj}(\cdot))} \quad \underbrace{\hspace{10em}}_{\text{gate_proj}(\cdot)}
 \end{array}$$

$$\begin{array}{c}
 \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \odot \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{0}^D \rightarrow \end{array} \right) \xrightarrow{\text{out_proj}(\cdot)} \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \odot \mathbf{x}_0 \rightarrow \\ \leftarrow (\mathbf{x}_0^T \mathbf{a}_i - \mathbf{b}_i) \mathbf{1}^D \rightarrow \end{array} \right)
 \end{array}$$

2452 Layer 2:

$$\begin{array}{l}
 2453 \\
 2454 \\
 2455 \\
 2456 \\
 2457 \\
 2458 \\
 2459 \\
 2460 \\
 2461 \\
 2462 \\
 2463 \\
 2464 \\
 2465 \\
 2466 \\
 2467 \\
 2468 \\
 2469 \\
 2470 \\
 2471 \\
 2472 \\
 2473
 \end{array}$$

$$\begin{array}{c}
 \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \odot \mathbf{x}_0 \rightarrow \\ \leftarrow (\mathbf{x}_0^T \mathbf{a}_i - \mathbf{b}_i) \mathbf{1}^D \rightarrow \end{array} \right) \odot \left(\begin{array}{c} \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1} \rightarrow \\ \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{a}_i \rightarrow \end{array} \right) = \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \odot \mathbf{x}_0 \rightarrow \\ \leftarrow (\mathbf{x}_0^T \mathbf{a}_i - \mathbf{b}_i) \mathbf{a}_i \rightarrow \end{array} \right) \\
 \underbrace{\hspace{10em}}_{\text{conv}(\text{in_proj}(\cdot))} \quad \underbrace{\hspace{10em}}_{\text{gate_proj}(\cdot)}
 \end{array}$$

$$\begin{array}{c}
 \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \odot \mathbf{x}_0 \rightarrow \\ \leftarrow (\mathbf{x}_0^T \mathbf{a}_i - \mathbf{b}_i) \mathbf{a}_i \rightarrow \end{array} \right) \xrightarrow{\text{out_proj}(\cdot)=\text{Identity}} \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \odot \mathbf{x}_0 \rightarrow \\ \leftarrow (\mathbf{x}_0^T \mathbf{a}_i - \mathbf{b}_i) \mathbf{a}_i \rightarrow \end{array} \right)
 \end{array}$$

2474 Layer 3:

$$\begin{array}{l}
 2475 \\
 2476 \\
 2477 \\
 2478 \\
 2479 \\
 2480 \\
 2481 \\
 2482 \\
 2483
 \end{array}$$

$$\begin{array}{c}
 \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \odot \mathbf{x}_0 \rightarrow \\ \leftarrow (\mathbf{x}_0^T \mathbf{a}_i - \mathbf{b}_i) \mathbf{a}_i \rightarrow \end{array} \right) \xrightarrow{\text{conv}(\text{in_proj}(\cdot))} \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \odot \mathbf{x}_0 \rightarrow \\ \leftarrow \sum_{i=1}^N (\mathbf{x}_0^T \mathbf{a}_i - \mathbf{b}_i) \mathbf{a}_i \rightarrow \end{array} \right)
 \end{array}$$

$$\begin{array}{l}
2484 \\
2485 \\
2486 \\
2487 \\
2488 \\
2489 \\
2490 \\
2491 \\
2492 \\
2493 \\
2494 \\
2495 \\
2496 \\
2497
\end{array}
\left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 \rightarrow \\ \leftarrow \mathbf{a}_i \odot \mathbf{x}_0 \rightarrow \\ \leftarrow \underbrace{\sum_{i=1}^N (\mathbf{x}_0^T \mathbf{a}_i - \mathbf{b}_i) \mathbf{a}_i}_{=\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_0)} \rightarrow \end{array} \right) \xrightarrow{\text{out_proj}(\cdot)} \left(\begin{array}{c} \leftarrow \mathbf{a}_i \rightarrow \\ \leftarrow \mathbf{b}_i \rightarrow \\ \leftarrow \mathbf{x}_0 - \eta \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{w}_0) \rightarrow \\ \leftarrow \mathbf{0}^D \rightarrow \\ \leftarrow \mathbf{0}^D \rightarrow \end{array} \right)$$

2498 **Causal BASECONV** This weight construction uses Equation 47 to compute the gradient descent
2499 update.

2500 We start with input:

$$\begin{array}{l}
2501 \\
2502 \\
2503 \\
2504 \\
2505 \\
2506 \\
2507 \\
2508
\end{array}
\mathbf{b} \equiv \begin{pmatrix} \mathbf{a}_1 & \dots & \mathbf{a}_N & \mathbf{0}^D \\ \mathbf{b}_1 & \dots & \mathbf{b}_N & 0 \\ \mathbf{0}^D & \dots & \mathbf{0}^D & \mathbf{x}_0 \end{pmatrix}$$

2509 We use two BASECONV layers to construct an initial embedding, after which each gradient descent
2510 update step will only require a single BASECONV layer.

2511 In the following construction, we use *flt* to denote the *flatten* operation, which maps an $M \times N$
2512 matrix to a MN -entry vector with the same elements.

2513 Layer 1:

$$\begin{array}{l}
2514 \\
2515 \\
2516 \\
2517 \\
2518 \\
2519 \\
2520 \\
2521 \\
2522 \\
2523 \\
2524 \\
2525 \\
2526 \\
2527 \\
2528 \\
2529 \\
2530 \\
2531 \\
2532 \\
2533 \\
2534 \\
2535 \\
2536 \\
2537
\end{array}
\left(\begin{array}{c} \mathbf{a}_1 \quad \dots \quad \mathbf{a}_N \quad \mathbf{0}^D \\ \mathbf{b}_1 \quad \dots \quad \mathbf{b}_N \quad 0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{x}_0 \\ \mathbf{a}_1 \quad \dots \quad \mathbf{a}_N \quad \mathbf{0}^D \\ \underbrace{\text{flt}(\mathbf{a}_1 \mathbf{1}^D)^T \quad \dots \quad \text{flt}(\mathbf{a}_N \mathbf{1}^D)^T \quad \text{flt}(\mathbf{0}^D (\mathbf{0}^D)^T)}_{\text{conv}(\text{in_proj}(\cdot))} \end{array} \right) \odot \left(\begin{array}{c} \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1} \rightarrow \\ \leftarrow \mathbf{1}^D \rightarrow \\ \mathbf{b}_1 \mathbf{1}^D \quad \dots \quad \mathbf{b}_N \mathbf{1}^D \quad \mathbf{0}^D \\ \underbrace{\text{flt}(\mathbf{1}^D \mathbf{a}_1^T) \quad \dots \quad \text{flt}(\mathbf{1}^D \mathbf{a}_N^T) \quad \text{flt}(\mathbf{0}^D (\mathbf{0}^D)^T)}_{\text{gate_proj}(\cdot)} \end{array} \right) =$$

$$\left(\begin{array}{c} \mathbf{a}_1 \quad \dots \quad \mathbf{a}_N \quad \mathbf{0}^D \\ \mathbf{b}_1 \quad \dots \quad \mathbf{b}_N \quad 0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{x}_0 \\ \mathbf{b}_1 \mathbf{a}_1 \quad \dots \quad \mathbf{b}_1 \mathbf{a}_N \quad \mathbf{0}^D \\ \text{flt}(\mathbf{a}_1 \mathbf{a}_1^T) \quad \dots \quad \text{flt}(\mathbf{a}_N \mathbf{a}_N^T) \quad \text{flt}(\mathbf{0}^D (\mathbf{0}^D)^T) \end{array} \right) \xrightarrow{\text{out_proj}=\text{Identity}} \left(\begin{array}{c} \mathbf{a}_1 \quad \dots \quad \mathbf{a}_N \quad \mathbf{0}^D \\ \mathbf{b}_1 \quad \dots \quad \mathbf{b}_N \quad 0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{x}_0 \\ \mathbf{b}_1 \mathbf{a}_1 \quad \dots \quad \mathbf{b}_1 \mathbf{a}_N \quad \mathbf{0}^D \\ \text{flt}(\mathbf{a}_1 \mathbf{a}_1^T) \quad \dots \quad \text{flt}(\mathbf{a}_N \mathbf{a}_N^T) \quad \text{flt}(\mathbf{0}^D (\mathbf{0}^D)^T) \end{array} \right)$$

2538 Layer 2:

$$\begin{array}{l}
 2540 \\
 2541 \\
 2542 \\
 2543 \\
 2544 \\
 2545 \\
 2546 \\
 2547 \\
 2548 \\
 2549 \\
 2550 \\
 2551 \\
 2552 \\
 2553 \\
 2554 \\
 2555 \\
 2556 \\
 2557 \\
 2558 \\
 2559 \\
 2560 \\
 2561
 \end{array}$$

$$\begin{array}{c}
 \left(\begin{array}{c} \mathbf{a}_1 \quad \dots \quad \mathbf{a}_N \quad \mathbf{0}^D \\ \mathbf{b}_1 \quad \dots \quad \mathbf{b}_N \quad 0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{x}_0 \\ \leftarrow \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \rightarrow \\ \leftarrow \sum_{i=1}^N \text{flt}(\mathbf{a}_i \mathbf{a}_i^T) \rightarrow \end{array} \right) \odot \underbrace{\left(\begin{array}{c} \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1} \rightarrow \\ \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1}^{D^2} \rightarrow \end{array} \right)}_{\text{gate_proj}(\cdot)} = \left(\begin{array}{c} \mathbf{a}_1 \quad \dots \quad \mathbf{a}_N \quad \mathbf{0}^D \\ \mathbf{b}_1 \quad \dots \quad \mathbf{b}_N \quad 0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{x}_0 \\ \leftarrow \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \rightarrow \\ \leftarrow \sum_{i=1}^N \text{flt}(\mathbf{a}_i \mathbf{a}_i^T) \rightarrow \end{array} \right) \\
 \underbrace{\hspace{10em}}_{\text{conv}(\text{in_proj}(\cdot))} \hspace{10em} \underbrace{\hspace{10em}}_{\text{gate_proj}(\cdot)} \hspace{10em} \underbrace{\hspace{10em}}_{\text{conv}(\text{in_proj}(\cdot))}
 \end{array}$$

$$\begin{array}{c}
 2552 \\
 2553 \\
 2554 \\
 2555 \\
 2556 \\
 2557 \\
 2558 \\
 2559 \\
 2560 \\
 2561
 \end{array}$$

$$\begin{array}{c}
 \left(\begin{array}{c} \mathbf{a}_1 \quad \dots \quad \mathbf{a}_N \quad \mathbf{0}^D \\ \mathbf{b}_1 \quad \dots \quad \mathbf{b}_N \quad 0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{x}_0 \\ \leftarrow \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \rightarrow \\ \leftarrow \sum_{i=1}^N \text{flt}(\mathbf{a}_i \mathbf{a}_i^T) \rightarrow \end{array} \right) \xrightarrow{\text{out_proj=Identity}} \left(\begin{array}{c} \mathbf{a}_1 \quad \dots \quad \mathbf{a}_N \quad \mathbf{0}^D \\ \mathbf{b}_1 \quad \dots \quad \mathbf{b}_N \quad 0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{x}_0 \\ \leftarrow \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \rightarrow \\ \leftarrow \sum_{i=1}^N \text{flt}(\mathbf{a}_i \mathbf{a}_i^T) \rightarrow \end{array} \right)
 \end{array}$$

2562 Now, we use a single BASECONV layer to implement a gradient descent update.

$$\begin{array}{l}
 2562 \\
 2563 \\
 2564 \\
 2565 \\
 2566 \\
 2567 \\
 2568 \\
 2569 \\
 2570 \\
 2571 \\
 2572 \\
 2573 \\
 2574 \\
 2575 \\
 2576 \\
 2577 \\
 2578 \\
 2579 \\
 2580 \\
 2581 \\
 2582
 \end{array}$$

$$\begin{array}{c}
 \left(\begin{array}{c} \mathbf{a}_1 \quad \dots \quad \mathbf{a}_N \quad \mathbf{0}^D \\ \mathbf{b}_1 \quad \dots \quad \mathbf{b}_N \quad 0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{x}_0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{1}^D \\ \leftarrow \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \rightarrow \\ \leftarrow \sum_{i=1}^N \text{flt}(\mathbf{a}_i \mathbf{a}_i^T) \rightarrow \\ \leftarrow \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \rightarrow \\ \leftarrow \sum_{i=1}^N \text{flt}(\mathbf{a}_i \mathbf{a}_i^T) \rightarrow \end{array} \right) \odot \underbrace{\left(\begin{array}{c} \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1} \rightarrow \\ \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1}^D \rightarrow \\ \leftarrow \mathbf{1}^{D^2} \rightarrow \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{1}^D \\ \mathbf{0}^{D^2} \quad \dots \quad \mathbf{0}^{D^2} \quad \text{flt}(\mathbf{1}^D \mathbf{x}_0^T) \end{array} \right)}_{\text{gate_proj}(\cdot)} = \left(\begin{array}{c} \mathbf{a}_1 \quad \dots \quad \mathbf{a}_N \quad \mathbf{0}^D \\ \mathbf{b}_1 \quad \dots \quad \mathbf{b}_N \quad 0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{x}_0 \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \mathbf{1}^D \\ \leftarrow \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \rightarrow \\ \leftarrow \sum_{i=1}^N \text{flt}(\mathbf{a}_i \mathbf{a}_i^T) \rightarrow \\ \mathbf{0}^D \quad \dots \quad \mathbf{0}^D \quad \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \\ \mathbf{0}^{D^2} \quad \dots \quad \mathbf{0}^{D^2} \quad \sum_{i=1}^N \text{flt}(\mathbf{a}_i (\mathbf{a}_i \odot \mathbf{x}_0)^T) \end{array} \right) \\
 \underbrace{\hspace{10em}}_{\text{conv}(\text{in_proj}(\cdot))} \hspace{10em} \underbrace{\hspace{10em}}_{\text{gate_proj}(\cdot)} \hspace{10em} \underbrace{\hspace{10em}}_{\text{conv}(\text{in_proj}(\cdot))}
 \end{array}$$

2583 Note that the gradient

$$\begin{array}{l}
 2584 \\
 2585 \\
 2586 \\
 2587
 \end{array}$$

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_0) = \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i - \left(\sum_{i=1}^N \mathbf{a}_i \mathbf{a}_i^T \right) \mathbf{w}_0$$

2588 can be written as a linear combination of the vector

$$\begin{array}{l}
 2589 \\
 2590 \\
 2591
 \end{array}$$

$$\left(\begin{array}{c} \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \\ \sum_{i=1}^N \text{flt}(\mathbf{a}_i (\mathbf{a}_i \odot \mathbf{x}_0)^T) \end{array} \right)$$

so we can write a weight construction for out_proj that updates $w_0 \rightarrow w_0 - \eta \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_0)$:

$$\begin{array}{c}
 2592 \\
 2593 \\
 2594 \\
 2595 \\
 2596 \\
 2597 \\
 2598 \\
 2599 \\
 2600 \\
 2601 \\
 2602 \\
 2603 \\
 2604 \\
 2605 \\
 2606 \\
 2607 \\
 2608 \\
 2609
 \end{array}
 \left(\begin{array}{cccc}
 \mathbf{a}_1 & \dots & \mathbf{a}_N & \mathbf{0}^D \\
 \mathbf{b}_1 & \dots & \mathbf{b}_N & 0 \\
 \mathbf{0}^D & \dots & \mathbf{0}^D & \mathbf{x}_0 \\
 \mathbf{0}^D & \dots & \mathbf{0}^D & \mathbf{1}^D \\
 \leftarrow \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \rightarrow \\
 \leftarrow \sum_{i=1}^N fll(\mathbf{a}_i \mathbf{a}_i^T) \rightarrow \\
 \mathbf{0}^D & \dots & \mathbf{0}^D & \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \\
 \mathbf{0}^{D^2} & \dots & \mathbf{0}^{D^2} & \sum_{i=1}^N fll(\mathbf{a}_i (\mathbf{a}_i \odot \mathbf{x}_0)^T)
 \end{array} \right)
 \xrightarrow[\text{out_proj}]{}
 \left(\begin{array}{cccc}
 \mathbf{a}_1 & \dots & \mathbf{a}_N & \mathbf{0}^D \\
 \mathbf{b}_1 & \dots & \mathbf{b}_N & 0 \\
 \mathbf{0}^D & \dots & \mathbf{0}^D & \mathbf{x}_0 - \eta \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_0) \\
 \mathbf{0}^D & \dots & \mathbf{0}^D & \mathbf{1}^D \\
 \leftarrow \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \rightarrow \\
 \leftarrow \sum_{i=1}^N fll(\mathbf{a}_i \mathbf{a}_i^T) \rightarrow \\
 \mathbf{0}^D & \dots & \mathbf{0}^D & \sum_{i=1}^N \mathbf{b}_i \mathbf{a}_i \\
 \mathbf{0}^{D^2} & \dots & \mathbf{0}^{D^2} & \sum_{i=1}^N fll(\mathbf{a}_i (\mathbf{a}_i \odot \mathbf{x}_0)^T)
 \end{array} \right)$$

2610 D.3.2 LOWER BOUNDS: BASECONV CONSTRUCTIONS ARE ASYMPTOTICALLY OPTIMAL

2611
 2612 Note that the non-causal weight construction in Appendix D.3.1 requires $O(1)$ layers and $O(D)$ state
 2613 size, while the causal weight construction in Appendix D.3.1 requires $O(1)$ layers and $O(D^2)$ state
 2614 size. Clearly the $O(D)$ state size requirement for non-causal models is tight, since one needs to store
 2615 the gradient $\nabla_{\mathbf{x}} \mathcal{L} \in \mathbb{R}^D$. In this section, we prove that the $O(D^2)$ state size requirement for causal
 2616 models is also asymptotically tight.

2617 **Theorem D.33.** *Any single-pass (causal) algorithm computing the gradient*

$$2618 \\
 2619 \\
 2620 \\
 2621 \\
 2622 \\
 2623 \\
 2624 \\
 2625 \\
 2626 \\
 2627 \\
 2628 \\
 2629 \\
 2630 \\
 2631 \\
 2632 \\
 2633 \\
 2634 \\
 2635 \\
 2636 \\
 2637 \\
 2638 \\
 2639 \\
 2640 \\
 2641 \\
 2642 \\
 2643 \\
 2644 \\
 2645$$

$$\nabla_{\mathbf{x}} \mathcal{L} = \sum_{j=1}^N b_j \mathbf{a}_j - \left(\sum_{j=1}^N \mathbf{a}_j \mathbf{a}_j^T \right) \mathbf{x}$$

2623 *given inputs $\{(\mathbf{a}_1, b_1), \dots, (\mathbf{a}_N, b_N)\}; \mathbf{x}$, with $(\mathbf{a}_i, b_i) \in \mathbb{R}^{(D+1)N}$ and $\mathbf{x} \in \mathbb{R}^D$, requires $\Omega(D^2)$*
 2624 *state size in the worst case, where $b_j \in \mathbb{R}$ and $\mathbf{a}_j, \mathbf{x} \in \mathbb{R}^D$.*

2626 *Proof.* For simplicity, we pick $N = D$ for large enough D .

2627 Since we can compute $\sum_{j=1}^D b_j \mathbf{a}_j$ in $O(D)$ space, we focus on computing the expensive
 2628 $\left(\sum_{j=1}^N \mathbf{a}_j \mathbf{a}_j^T \right) \mathbf{x}$ term. Assume there exists a single-pass algorithm \mathcal{A} that computes
 2629 $\left(\sum_{j=1}^N \mathbf{a}_j \mathbf{a}_j^T \right) \mathbf{x}$ exactly for all choices of $\mathbf{a}_1, \dots, \mathbf{a}_D, \mathbf{x} \in \mathbb{R}^D$. Now consider the following
 2630 two claims:
 2631
 2632

- 2633 1. Define s_D to be the state of the algorithm after seeing $\mathbf{a}_1, \dots, \mathbf{a}_D$. Then we claim that s_D
 2634 must have enough information to exactly reconstruct $\mathbf{M}_D := \sum_{j=1}^D \mathbf{a}_j \mathbf{a}_j^T$.

2635 This follows since the algorithm must be correct for any value $\mathbf{x} \in \mathbb{R}^D$ takes on. In
 2636 particular, setting $\mathbf{x} = \mathbf{e}_i$ for $i \in [D]$, we observe that the algorithm must be able to exactly
 2637 recover $\mathbf{M}_D \mathbf{e}_i = \mathbf{M}_D[:, i]$, $i \in [D]$.

- 2640 2. The space of matrices

$$2641 \\
 2642 \\
 2643 \\
 2644 \\
 2645$$

$$\left\{ \sum_{j=1}^D \mathbf{a}_j \mathbf{a}_j^T \right\}$$

over all choices of $\mathbf{a}_j \in \mathbb{R}^D$, $j \in [d]$ contains the set of all real symmetric matrices in
 $\mathbb{R}^{D \times D}$.

2646 This holds since for any real symmetric matrix A , we can obtain a set of possible \mathbf{a}_j 's via
 2647 its eigendecomposition [Strang \(2012\)](#):
 2648

$$2649 \quad A = Q\Lambda Q^T = \sum_{j=1}^D \mathbf{a}_j \mathbf{a}_j^T$$

2650
 2651 where $\mathbf{a}_j = \sqrt{\lambda_j} Q[:, j]$.
 2652
 2653

2654 From the first claim, we conclude that s_D must contain enough information to be able to recover M_D
 2655 for any possible value M_D can take on (over all choices of $\mathbf{a}_1, \dots, \mathbf{a}_D \in \mathbb{R}^D$). From the second
 2656 claim, we have that the space of possible values of M_D includes the set of all possible real symmetric
 2657 matrices. Since we know that this set requires $\frac{(D)(D+1)}{2}$ parameters to represent, we can conclude
 2658 that $|s_D| \geq \frac{(D)(D+1)}{2} \geq \Omega(D^2)$. □
 2659
 2660
 2661
 2662
 2663
 2664
 2665
 2666
 2667
 2668
 2669
 2670
 2671
 2672
 2673
 2674
 2675
 2676
 2677
 2678
 2679
 2680
 2681
 2682
 2683
 2684
 2685
 2686
 2687
 2688
 2689
 2690
 2691
 2692
 2693
 2694
 2695
 2696
 2697
 2698
 2699

2700 D.4 BASECONV AND JACKSON’S THEOREM
2701

2702 In this section we prove BASECONV’s ability to approximate arbitrary univariate and multivariate
2703 smooth functions.

2704 We start with a special case of smooth functions that apply entry-wise univariate smooth functions:
2705

2706 **Definition D.34.** Let $\bar{f} : [-1, 1] \rightarrow \mathbb{R}$ be a (k, L) -smooth univariate function. Then define

$$2707 f : [-1, 1]^{N \times D} \rightarrow \mathbb{R}^{N \times D}$$

2708 as follows. For all $0 \leq i < N$, $0 \leq j < D$, and $\mathbf{u} \in [-1, 1]^{N \times D}$:

$$2709 (f(\mathbf{u}))[i, j] = \bar{f}(\mathbf{u}[i, j]).$$

2710 Now we will state a simple observation on BASECONV’s ability to approximate these functions.
2711

2712 **Lemma D.35.** For any smooth function f as defined in [Definition D.34](#), let $g(\mathbf{x}) = P_{\bar{f}}(\mathbf{x})$ with $P_{\bar{f}}$
2713 being the polynomial from [Corollary D.12](#). Then for all $\mathbf{x} \in [-1, 1]^{N \times D}$,

$$2714 \|g(\mathbf{x}) - f(\mathbf{x})\|_{\infty} \leq \epsilon.$$

2715 *Proof.* Follows from [Definitions D.7](#) and [D.34](#) and [Corollary D.12](#). □

2716 Next we will state a construction of an arithmetic circuit for a function that applies a univariate
2717 polynomial to all entries in $[-1, 1]^{N \times D}$:

2718 **Lemma D.36.** Let $P(X)$ be a degree d univariate polynomial. Then there is a
2719 $(ND, O(ND), O(d), ND)$ -circuit to compute $P(\mathbf{u})$ where $P(\mathbf{u})$ is defined as follows. For an
2720 input $\mathbf{u} \in [-1, 1]^{N \times D}$,

$$2721 P(\mathbf{u})[i, j] = P(\mathbf{u}[i, j]).$$

2722 *Proof.* Let the univariate polynomial be

$$2723 P(X) = \sum_{i=0}^d c_i X^i$$

2724 where coefficients $c_i \in \mathbb{R}$.

2725 Next we state the natural arithmetic circuit to compute $P(x)$ for $x \in \mathbb{R}$ in [Algorithm 2](#):

2726 **Algorithm 2** circuit $\mathcal{C}_P(x)$:

2727 1: $s_0 \leftarrow c_0$
2728 2: $m_0 \leftarrow 1$
2729 3: **for** $j = 1, 2, \dots, d$ **do**
2730 4: $m_j \leftarrow m_{j-1} \cdot x$ ▷ Multiplication gate
2731 5: $t_j \leftarrow c_j \cdot m_j$ ▷ Multiplication gate
2732 6: $s_j \leftarrow s_{j-1} + t_j$ ▷ Addition gate
2733 7: **return** s_d ▷ s_d is the output gate

2734 Next we apply the above circuit in parallel to form the circuit that computes $P(\mathbf{u})$ in [Algorithm 3](#):

2735 **Algorithm 3** Circuit for $P(\mathbf{u})$:

2736 1: **for** $i = 0, 1, \dots, N - 1$ **do**
2737 2: **for** $j = 0, 1, \dots, D - 1$ **do**
2738 3: $\mathbf{z}[i, j] = \mathcal{C}_P(\mathbf{u}[i, j])$ ▷ Do this in parallel
2739 4: **return** \mathbf{z} ▷ \mathbf{z} is the output matrix

Looking at [Algorithm 2](#), the depth of the circuit is $3d$, or $O(d)$, since that is the bound on iterations of the for loop, and each iteration we compute 3 sequential operations. Therefore it's a $(1, O(d), O(d), O(1))$ -circuit.

For [Algorithm 3](#), The width is $O(ND)$, since we have our input of size $N \times D$, which goes through the circuit in parallel, as stated in [Algorithm 3](#). Therefore we have an $(ND, O(ND), O(d), O(ND))$ -circuit that computes $P(\mathbf{u})$. \square

Since BASECONV has the ability to represent any arithmetic circuit, we get the following:

Corollary D.37. *We can implement $P(\mathbf{u})$ (where $P(\mathbf{u})$ is as defined in [Lemma D.36](#)) when $\deg(P) = d$ with a $(N, O(d \log(ND)), D, O(ND), D)$ – BASECONV.*

Proof. Follows from [Lemma D.36](#) giving us the $(ND, O(ND), O(d), O(ND))$ -circuit for an arbitrary polynomial and [Theorem D.30](#) gives us the BASECONV model to implement the circuit. \square

We will prove a tighter bound showing we can represent $P(\mathbf{u})$ using a constant number of BASECONV layers (for constant $\deg(P)$):

Theorem D.38. *We can implement $P(\mathbf{u})$ when $\deg(P) = d$ with an $(O(N), O(d), D, O(N), D)$ – BASECONV model.*

Proof. We will convert the steps done in [Algorithm 2](#) to layers of BASECONV. Since [Algorithm 3](#) is essentially running [Algorithm 2](#) in parallel over all entries of input $\mathbf{u} \in [-1, 1]^{N \times D}$, the latter happens automatically in our BASECONV implementation.

For this proof, define

$$P_j(X) = X^j$$

and let C_i be the matrix of size $N \times D$ and all the entries are c_i .

We expand the input to our BASECONV layers as follows,

$$\mathbf{u} = \begin{pmatrix} \mathbf{u}' \\ \mathbf{0}^{3N \times D} \end{pmatrix}.$$

This means that the size of the internal dimension of our BASECONV layers will be $(4N, D)$.

To begin iterations of the for loop we need to store initial values into the extra space in \mathbf{u} . Taking us from

$$\mathbf{u} = \begin{pmatrix} \mathbf{u}' \\ \mathbf{0}^{N \times D} \\ \mathbf{0}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix} \rightarrow \begin{pmatrix} \mathbf{u} \\ \mathbf{1}^{N \times D} \\ \mathbf{1}^{N \times D} \\ \mathbf{C}_0 \end{pmatrix} =: \mathbf{u}_0$$

We do this via $\text{BASECONV}(\mathbf{u}', \mathbf{I}^{D \times D}, \begin{pmatrix} \mathbf{0}^{N \times D} \\ \mathbf{1}^{N \times D} \\ \mathbf{1}^{N \times D} \\ \mathbf{C}_0 \end{pmatrix}, \mathbf{0}^{4N \times D}, \mathbf{1}^{4N \times D})$ which computes

$$\left(\begin{pmatrix} \mathbf{u} \\ \mathbf{0}^{N \times D} \\ \mathbf{0}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix} \mathbf{I}^{D \times D} + \begin{pmatrix} \mathbf{0}^{N \times D} \\ \mathbf{1}^{N \times D} \\ \mathbf{1}^{N \times D} \\ \mathbf{C}_0 \end{pmatrix} \right) \odot \left(\mathbf{0}^{4N \times D} * \begin{pmatrix} \mathbf{u} \\ \mathbf{0}^{N \times D} \\ \mathbf{0}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix} + \mathbf{1}^{4N \times D} \right).$$

The above simplifies to

$$\left(\begin{pmatrix} \mathbf{u} \\ \mathbf{0}^{N \times D} \\ \mathbf{0}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix} + \begin{pmatrix} \mathbf{0}^{N \times D} \\ \mathbf{1}^{N \times D} \\ \mathbf{1}^{N \times D} \\ \mathbf{C}_0 \end{pmatrix} \right) \odot (\mathbf{1}^{4N \times D}),$$

which gives us

$$\begin{pmatrix} \mathbf{u} \\ \mathbf{1}^{N \times D} \\ \mathbf{1}^{N \times D} \\ \mathbf{C}_0 \end{pmatrix} =: \mathbf{u}_0,$$

2808 as desired
 2809

2810 This was done with a $(4N, 1, D, 4N, D)$ – BASECONV layer.

2811 Our goal is, at the end of iteration j to compute $\mathbf{u}_j \in \mathbb{R}^{4N \times D}$ such that,

$$2812 \mathbf{u}_j = \begin{pmatrix} \mathbf{u} \\ P_j(\mathbf{u}) \\ C_j \odot P_j(\mathbf{u}) \\ C_0 + C_1 \odot P_1(\mathbf{u}) + \cdots + C_j \odot P_j(\mathbf{u}) \end{pmatrix}.$$

2817 We will view the above matrix in terms of the variables in the [Algorithm 2](#) as follows

$$2818 \begin{pmatrix} \mathbf{u} \\ P_j(\mathbf{u}) \\ C_j \odot P_j(\mathbf{u}) \\ C_0 + C_1 \odot P_1(\mathbf{u}) + \cdots + C_j \odot P_j(\mathbf{u}) \end{pmatrix} =: \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_j \\ \mathbf{t}_j \\ \mathbf{s}_j \end{pmatrix}.$$

2824 The for loop runs for values of $1 \leq j \leq d$ which the remainder of this proof will replicate. There
 2825 are three lines in the for loop in [Algorithm 2](#) which we will cover how these operations happen in
 2826 constant number of BASECONV layers.

2827 In line 4, the first line in the for loop computes

$$2828 \mathbf{u}_{j-1} = \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_{j-1} \\ \mathbf{t}_{j-1} \\ \mathbf{s}_{j-1} \end{pmatrix} \rightarrow \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_j \\ \mathbf{t}_{j-1} \\ \mathbf{s}_{j-1} \end{pmatrix} =: \mathbf{u}_j^{(1)}.$$

2833 Note that $\mathbf{m}_j = \mathbf{m}_{j-1} \odot \mathbf{u}$.

2835 We use the remember primitive to compute $\mathbf{u}_j^{(1)}$ from \mathbf{u}_{j-1} . Define $f : \mathbb{R}^{2N \times D} \rightarrow \mathbb{R}^{2N \times D}$ as
 2836 follows

$$2837 f \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_{j-1} \end{pmatrix} = \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_{j-1} \odot \mathbf{u} \end{pmatrix}.$$

2839 If we can compute f with BASECONV layers then we can compute $\mathbf{u}_j^{(1)}$ for \mathbf{u}_{j-1} by calling
 2840 remember($\mathbf{u}_j, 0, 2N - 1, f$).

2842 We show BASECONV $\left(\begin{pmatrix} \mathbf{u} \\ \mathbf{m}_j \end{pmatrix}, \mathbf{I}^{D \times D}, \mathbf{0}^{2N \times D}, \mathbf{H}, \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix} \right)$ maps

$$2843 \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_{j-1} \end{pmatrix} \rightarrow \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_j \end{pmatrix},$$

2847 where \mathbf{H} is defined as in [Proposition D.24](#). We plug the matrices into the BASECONV layer as
 2848 follows:

$$2849 \left(\begin{pmatrix} \mathbf{u} \\ \mathbf{m}_{j-1} \end{pmatrix} \cdot \mathbf{I}^{D \times D} + \mathbf{0}^{2N \times D} \right) \odot \left(\mathbf{H} * \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_{j-1} \end{pmatrix} + \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix} \right).$$

2851 We know from [Proposition D.24](#) that this convolution operation is a shift down by N rows. Therefore
 2852 the above simplifies to

$$2853 \left(\begin{pmatrix} \mathbf{u} \\ \mathbf{m}_{j-1} \end{pmatrix} \cdot \mathbf{I}^{D \times D} + \mathbf{0}^{2N \times D} \right) \odot \left(\begin{pmatrix} \mathbf{0}^{N \times D} \\ \mathbf{u} \end{pmatrix} + \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix} \right),$$

2856 which simplifies to

$$2857 \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_{j-1} \end{pmatrix} \odot \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{u} \end{pmatrix} = \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_{j-1} \odot \mathbf{u} \end{pmatrix} = f \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_j \end{pmatrix},$$

2858 as desired. Therefore by [Proposition D.26](#), line 4 can be computed by $(4N, 8, D, 4N, D)$ –
 2860 BASECONV.
 2861

2862 For line 5 of the for loop we need to compute
 2863

$$2864 \mathbf{u}_j^{(1)} = \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_j \\ \mathbf{t}_{j-1} \\ \mathbf{s}_{j-1} \end{pmatrix} \rightarrow \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_j \\ \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix} =: \mathbf{u}_j^{(2)}.$$

2865
 2866
 2867 Note that $\mathbf{t}_j = \mathbf{C}_j \odot \mathbf{m}_j$.

2868
 2869 To do this we will use three BASECONV layers. We use the `remember` primitive to compute $\mathbf{u}_j^{(2)}$
 2870 from $\mathbf{u}_j^{(1)}$. Define $g : \mathbb{R}^{2N \times D} \rightarrow \mathbb{R}^{2N \times D}$ as follows,
 2871

$$2872 g \begin{pmatrix} \mathbf{m}_j \\ \mathbf{t}_{j-1} \end{pmatrix} = \begin{pmatrix} \mathbf{m}_j \\ \mathbf{C}_j \odot \mathbf{m}_j \end{pmatrix}.$$

2873
 2874 If we can compute g with BASECONV layers then we can compute $\mathbf{u}_j^{(2)}$ for \mathbf{u}_{j-1} by calling
 2875 `remember`($\mathbf{u}_j^{(1)}$, N , $3N - 1$, g).
 2876
 2877

2878 Indeed, we show the g can be computed by first computing
 2879 `BASECONV`($\begin{pmatrix} \mathbf{m}_j \\ \mathbf{t}_{j-1} \end{pmatrix}$, $\mathbf{I}^{D \times D}$, $\mathbf{0}^{2N \times D}$, $\mathbf{0}^{2N \times D}$, $\begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix}$):
 2880

$$2881 \left(\begin{pmatrix} \mathbf{m}_j \\ \mathbf{t}_{j-1} \end{pmatrix} \cdot \mathbf{I}^{D \times D} + \mathbf{0}^{2N \times D} \right) \odot \left(\mathbf{0}^{2N \times D} * \begin{pmatrix} \mathbf{m}_j \\ \mathbf{t}_{j-1} \end{pmatrix} + \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix} \right),$$

2882
 2883 which simplifies to

$$2884 \left(\begin{pmatrix} \mathbf{m}_j \\ \mathbf{t}_{j-1} \end{pmatrix} \right) \odot \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix}.$$

2885
 2886 This results in

$$2887 \begin{pmatrix} \mathbf{m}_j \\ \mathbf{0}^{N \times D} \end{pmatrix}.$$

2888
 2889 We pass into the next layer, `BASECONV`($\begin{pmatrix} \mathbf{m}_j \\ \mathbf{0}^{N \times D} \end{pmatrix}$, $\mathbf{I}^{D \times D}$, $\begin{pmatrix} \mathbf{0}^{N \times D} \\ \mathbf{1}^{N \times D} \end{pmatrix}$, \mathbf{H} , $\begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix}$) where \mathbf{H} is
 2890 defined as in [Proposition D.24](#):
 2891

$$2892 \left(\begin{pmatrix} \mathbf{m}_j \\ \mathbf{0}^{N \times D} \end{pmatrix} \cdot \mathbf{I}^{D \times D} + \begin{pmatrix} \mathbf{0}^{N \times D} \\ \mathbf{1}^{N \times D} \end{pmatrix} \right) \odot \left(\mathbf{H} * \begin{pmatrix} \mathbf{m}_j \\ \mathbf{0}^{N \times D} \end{pmatrix} + \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix} \right).$$

2893
 2894 Since the kernel \mathbf{H} is as in [Proposition D.24](#), this simplifies to

$$2895 \left(\begin{pmatrix} \mathbf{m}_j \\ \mathbf{1}^{N \times D} \end{pmatrix} \odot \left(\begin{pmatrix} \mathbf{0}^{N \times D} \\ \mathbf{m}_j \end{pmatrix} + \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{0}^{N \times D} \end{pmatrix} \right) \right).$$

2896
 2897 The above simplifies further to

$$2898 \begin{pmatrix} \mathbf{m}_j \\ \mathbf{1}^{N \times D} \end{pmatrix} \odot \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{m}_j \end{pmatrix},$$

2900
 2901 which results in:

$$2902 \begin{pmatrix} \mathbf{m}_j \\ \mathbf{m}_j \end{pmatrix}.$$

2903
 2904 We pass the above to `BASECONV`($\begin{pmatrix} \mathbf{m}_j \\ \mathbf{m}_j \end{pmatrix}$, $\mathbf{I}^{D \times D}$, $\mathbf{0}^{2N \times D}$, $\mathbf{0}^{2N \times D}$, $\begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{C}_j \end{pmatrix}$):
 2905

$$2906 \left(\begin{pmatrix} \mathbf{m}_j \\ \mathbf{m}_j \end{pmatrix} \cdot \mathbf{I}^{D \times D} + \mathbf{0}^{2N \times D} \right) \odot \left(\mathbf{0}^{2N \times D} * \begin{pmatrix} \mathbf{m}_j \\ \mathbf{m}_j \end{pmatrix} + \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{C}_j \end{pmatrix} \right)$$

2907
 2908 which simplifies to

$$2909 \begin{pmatrix} \mathbf{m}_j \\ \mathbf{m}_j \end{pmatrix} \odot \begin{pmatrix} \mathbf{1}^{N \times D} \\ \mathbf{C}_j \end{pmatrix}.$$

2910
 2911 The above results in

$$2912 \begin{pmatrix} \mathbf{m}_j \\ \mathbf{C}_j \odot \mathbf{m}_j \end{pmatrix} = g \begin{pmatrix} \mathbf{m}_j \\ \mathbf{t}_{j-1} \end{pmatrix},$$

2913
 2914
 2915

2916 as desired.

2917 Therefore by [Corollary D.27](#), line 5 was computed by $(4N, O(1), D, 4N, D) - \text{BASECONV}$.

2918 For line 6, the final line of the for loop, we want

$$2919 \mathbf{u}_j^{(2)} = \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_j \\ \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix} \rightarrow \begin{pmatrix} \mathbf{u} \\ \mathbf{m}_j \\ \mathbf{t}_j \\ \mathbf{s}_j \end{pmatrix} =: \mathbf{u}_j.$$

2920 Note that $\mathbf{s}_j = \mathbf{s}_{j-1} + \mathbf{t}_j$

2921 Define function $h : \mathbb{R}^{2N \times D} \rightarrow \mathbb{R}^{2N \times D}$ as follows,

$$2922 h \begin{pmatrix} \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix} = \begin{pmatrix} \mathbf{t}_j \\ \mathbf{s}_{j-1} + \mathbf{t}_j \end{pmatrix}.$$

2923 If we can compute h with `BASECONV` layers then we can compute \mathbf{u}_j for \mathbf{u}_{j-1} by calling
2924 `remember($\mathbf{u}_j^{(2)}, 2N, 4N - 1, h$)`.

2925 Indeed we show that h can be computed by computing
2926 `BASECONV` $\left(\begin{pmatrix} \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix}, \mathbf{0}^{D \times D}, \mathbf{1}^{2N \times D}, \overline{\mathbf{H}}, \mathbf{0}^{2N \times D} \right)$, where kernel $\overline{\mathbf{H}} \in \mathbb{R}^{2N \times D}$ is defined
2927 as:

$$2928 \overline{\mathbf{H}}[k, :] \equiv \begin{cases} \mathbf{1}^D & \text{if } k \in \{0, N\} \\ \mathbf{0}^D & \text{otherwise.} \end{cases}$$

2929 This layer computes

$$2930 \left(\begin{pmatrix} \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix} \cdot \mathbf{0}^{2N \times D} + \mathbf{1}^{2N \times D} \right) \odot \left(\overline{\mathbf{H}} * \begin{pmatrix} \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix} + \mathbf{0}^{2N \times D} \right).$$

2931 This simplifies to

$$2932 (\mathbf{1}^{2N \times D}) \odot \left(\overline{\mathbf{H}} * \begin{pmatrix} \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix} \right) = \left(\overline{\mathbf{H}} * \begin{pmatrix} \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix} \right).$$

2933 Now we compute this convolution for column i , $0 \leq i < 2N$. For notational convenience, let
2934 $\begin{pmatrix} \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix}$ be noted as matrix \mathbf{V} . Then we have:

$$2935 \overline{\mathbf{H}}[:, i] * \mathbf{V}[:, i] = \text{coeff} \left((1 + X^N) \mathbf{V}[:, i](X) \pmod{X^{2N}} \right),$$

2936 where $(1 + X^N)$ is the polynomial representation of the columns of $\overline{\mathbf{H}}$ (since there's a one in the 0th
2937 index and a one in the N th index of each column).

2938 The expression simplifies to

$$2939 \text{coeff} \mathbf{V}[:, i](X) + \mathbf{V}[:, i](X) X^N \pmod{X^{2N}},$$

2940 which can be broken down to

$$2941 \text{coeff} \left((\mathbf{V}[0][i] + \mathbf{V}[1][i]X + \dots + \mathbf{V}[2N-1][i]X^{2N-1}) \pmod{X^{2N}} \right) \\ 2942 + \text{coeff} \left((\mathbf{V}[0][i]X^N + \mathbf{V}[1][i]X^{N+1} + \dots + \mathbf{V}[2N-1][i]X^{3N-1}) \pmod{X^{2N}} \right)$$

2943 with the lower order terms in the second coefficient vector being zeros,

$$2944 \text{coeff} \left((\mathbf{V}[0][i] + \mathbf{V}[1][i]X + \dots + \mathbf{V}[2N-1][i]X^{2N-1}) \pmod{X^{2N}} \right) \\ 2945 + \text{coeff} \left((0 + 0X + \dots + 0X^{N-1} + \mathbf{V}[0][i]X^N + \dots + \mathbf{V}[2N-1][i]X^{3N-1}) \pmod{X^{2N}} \right)$$

2946 After taking $\pmod{X^{2N}}$ we get

$$2947 \text{coeff} \left(\mathbf{V}[0][i] + \mathbf{V}[1][i]X + \dots + \mathbf{V}[2N-1][i]X^{2N-1} \right) \\ 2948 + \text{coeff} \left(0 + 0X + \dots + 0X^{N-1} \mathbf{V}[0][i]X^N + \dots + \mathbf{V}[N-1][i]X^{2N-1} \right)$$

The first set of coefficients is the input matrix as is. And the second one is the input matrix shifted down as seen in [Proposition D.24](#). Therefore when we add these vectors we are doing

$$\begin{pmatrix} \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix} + \begin{pmatrix} \mathbf{0}^{N \times D} \\ \mathbf{t}_j \end{pmatrix} = h \begin{pmatrix} \mathbf{t}_j \\ \mathbf{s}_{j-1} \end{pmatrix},$$

as desired. Therefore by [Proposition D.26](#), line 6 is computed with by $(4N, 1, D, 4N, D) - \text{BASECONV}$.

The \mathbf{s}_d matrix gives us $\mathbf{C}_0 + \mathbf{C}_1 \odot \mathbf{m}_1 + \dots + \mathbf{C}_d \odot \mathbf{m}_d$. Recalling that

$$\mathbf{C}_0 + \mathbf{C}_1 \odot \mathbf{m}_1 + \dots + \mathbf{C}_d \odot \mathbf{m}_d \equiv \sum_{j=0}^d \mathbf{C}_j \odot \mathbf{u}^j = P(\mathbf{u}),$$

and hence \mathbf{s}_d is our desired output.

We have d layers, each consisting of $O(1)$ BASECONV layers. Giving us $O(d)$ many layers to implement [Algorithm 2](#).

Therefore, via the ability to stack BASECONV layers to do function composition, the for loop was computed by a $(4N, O(d), D, 4N, D) - \text{BASECONV}$, as desired. \square

The following states BASECONV’s ability to approximate a univariate smooth function:

Proposition D.39. *Let f be the (k, L) -smooth function defined in [Definition D.34](#). Then there is a $(N, O(\sqrt[k]{\frac{L}{\epsilon}}) + k, D, (ND), D) - \text{BASECONV}$ model that approximates f within error ϵ .*

Proof. Follows from [Corollary D.12](#), [Lemma D.35](#), and [Theorem D.38](#). \square

D.5 MULTIVARIATE FUNCTION APPROXIMATION

We begin by defining more multivariate notation.

We consider the following multivariate functions:

Definition D.40. For $0 \leq i < N, 0 \leq j < D$, let $\bar{f}_{i,j} : [-1, 1]^{N \times D} \rightarrow \mathbb{R}$ be a (k, L) -smooth multivariate function. Then define

$$f(\mathbf{x}) : [-1, 1]^{N \times D} \rightarrow \mathbb{R}^{N \times D}$$

as follows. For all $0 \leq i < N, 0 \leq j < D, \mathbf{u} \in [-1, 1]^{N \times D}$ define

$$f(\mathbf{u})[i, j] := \bar{f}_{i,j}(\mathbf{u}).$$

Lemma D.41. *For any smooth function f as defined in [Definition D.40](#), let $g(X_1, \dots, X_{N \times D}) = P_{\bar{f}}(X_1, \dots, X_{N \times D})$ be the polynomial from [Corollary D.14](#). Then for all $\mathbf{x} \in [-1, 1]^{N \times D}$,*

$$\|g(\mathbf{x}) - f(\mathbf{x})\|_{\infty} \leq \epsilon.$$

Proof. Follows from [Definitions D.7](#) and [D.40](#) and [Corollary D.14](#). \square

Next we will state a construction for an arithmetic circuit for a function that takes a $[-1, 1]^{N \times D}$ variable input:

Lemma D.42. *Let $P(X)$ be a degree d multivariate polynomial. Then there is a $(n, O(d \cdot n^d), O(d \log(n)), O(n^d))$ -circuit to compute $P(\mathbf{u})$ on any input $\mathbf{u} \in [-1, 1]^n$.*

Proof. Let the multivariate polynomial be as defined in [Definition D.6](#). We build the circuit to compute this in [Algorithm 4](#),

3024
3025
3026
3027
3028
3029
3030
3031
3032
3033
3034
3035
3036
3037
3038
3039
3040
3041
3042
3043
3044
3045
3046
3047
3048
3049
3050
3051
3052
3053
3054
3055
3056
3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077

Algorithm 4 circuit $\mathcal{C}_P(\mathbf{x})$:

```

1: for  $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{Z}_{\geq 0}^n$  such that  $\sum_{i=1}^n \alpha_i \leq d$  do
2:    $m_\alpha \leftarrow 1$ 
3:   for  $i = 1, 2, \dots, n$  do ▷ Done in parallel
4:     if  $\alpha_i \neq 0$  then
5:        $m_\alpha \leftarrow m_\alpha \cdot x_i^{\alpha_i}$ 
6:    $t_\alpha \leftarrow c_\alpha \cdot m_\alpha$ 
7: for  $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{Z}_{\geq 0}^n$  such that  $\sum_{i=1}^n \alpha_i \leq d$  do
8:    $s \leftarrow \sum t_\alpha$  ▷ Done in parallel
9: return  $s$ 

```

We compute the for loop starting on line 3 by making multiplications in parallel. Therefore obtaining a depth of $O(\log(d))$. We also have the for loop starting on line 7, making pairwise addition operations, resulting in a depth of $O(d \log(n))$. \square

We again use the result that BASECONV can represent any arithmetic circuit to get:

Corollary D.43. *We can implement $P(\mathbf{u})$ (where $P(\mathbf{u})$ is as defined in Lemma D.42) when $\deg(P(X_1, \dots, X_{ND})) = d$ with a $(N, O(d \log(ND)), D, O((ND)^d), D) - \text{BASECONV}$ where $\mathbf{u} \in [-1, 1]^{N \times D}$.*

Proof. Lemma D.42 gives us the arithmetic circuit that computes this polynomial. Then via Theorem D.30 we get a $(N, O(d \log(ND)), D, O((ND)^d), D) - \text{BASECONV}$ model to implement the circuit. \square

Finally we state BASECONV's ability to approximate multivariate smooth functions:

Proposition D.44. *Let f be the function defined in Definition D.40. Then there is a $(N, O(d \log(ND)), D, O((ND)^d), D) - \text{BASECONV}$ model that approximates f to within error ϵ , with $d = O_k(\sqrt[k]{\frac{NDE}{\epsilon}})$.*

Proof. We get the existence of a polynomial that approximates f for some ϵ from Corollary D.14. Then via Corollary D.43 we get that we can represent any polynomial, implying $(N, O(d \log(ND)), D, O((ND)^d), D) - \text{BASECONV}$ represents any polynomial that approximates the multivariate smooth function f . \square

3078 D.6 NOTATION FOR D.7
3079

3080 This notation section is for the succeeding subsection which will prove we can recover the functions
3081 Square and Linear exactly given some assumptions and expected gradients of their respective loss
3082 functions being 0.

3083 When indexing an entry of a 2 dimensional matrix, \mathbf{A} we denote it as $A_{i,j}$ where i is the row number
3084 and j is the column.
3085

3086 The following are the parameters of a BASECONV layer.
3087

- 3088 1. We will use 0 indexing
- 3089 2. $[N] = \{0, 1, \dots, N - 1\}$
- 3090 3. $\mathbf{W} = \{\mathbf{W}_{i,j}\} \in \mathbb{R}^{d \times d}$
- 3091 4. $\mathbf{K} = \{\mathbf{K}_{i,j}\} \in \mathbb{R}^{N \times d}$
- 3092 5. $\mathbf{B}^{(1)} = \{\mathbf{B}_{i,j}^{(1)}\} \in \mathbb{R}^{N \times d}$
- 3093 6. $\mathbf{B}^{(2)} = \{\mathbf{B}_{i,j}^{(2)}\} \in \mathbb{R}^{N \times d}$
- 3094 7. $\boldsymbol{\theta} = (\mathbf{W}, \mathbf{K}, \mathbf{B}^{(1)}, \mathbf{B}^{(2)})$
- 3095 8. $\mathbf{u}_i \stackrel{\text{def}}{=} \mathbf{u}[i, :]$
- 3096 9. Input to a BASECONV layer, $\mathbf{u} = \{\mathbf{u}_{i,j}\} \in \mathbb{R}^{N \times d}$
- 3097 10. $\mathbb{E}[\{\mathbf{M}_{i,j}\}] = \{\mathbb{E}[\mathbf{M}_{i,j}]\}$
- 3098 11. BASECONV layer operation,

$$3100 \quad \mathbf{Z} = \text{BASECONV}(\mathbf{W}, \mathbf{K}, \mathbf{B}^{(1)}, \mathbf{B}^{(2)}, \mathbf{u}) \stackrel{\text{def}}{=} (\mathbf{u}\mathbf{W} + \mathbf{B}^{(1)}) \odot (\mathbf{K} * \mathbf{u} + \mathbf{B}^{(2)}) \quad (48)$$

- 3107 12. Our target function:

$$3108 \quad f : \mathbb{R}^{N \times d} \rightarrow \mathbb{R}^{N \times d}$$

3109 We'll denote $(f(\mathbf{u})) [i, j]$ by $f(\mathbf{u})_{i,j}$.

3110 We need to define the training input distribution to begin talking about expected values of the layers.
3111

3112 D.6.1 TRAINING INPUT DISTRIBUTION

- 3113 1. Let Δ be the training distribution on $\mathbb{R}^{N \times d}$

3114 **Assumption D.45.** Given a monomial $\mathbb{E}[\Pi_e(\mathbf{u}_{i_e, j_e})^{me}] = 0$ if me is odd. Otherwise,
3115 $\mathbb{E}[\Pi_e(\mathbf{u}_{i_e, j_e})^{me}] > 0$.

3116 **Assumption D.46.** Assume that the training data is generated as

- 3117 • $\mathbf{u} \sim \Delta$ as input
- 3118 • Output is $\mathbf{y} = f(\mathbf{u}) + \mathcal{E}$ where $\mathcal{E} = \{\mathcal{E}_{i,j}\} \in \mathbb{R}^{N \times d}$ is the random error matrix such that
 - 3119 – The distribution on \mathcal{E} and Δ are independent. (Call the distribution on \mathcal{E} to be $\Delta_{\mathcal{E}}$)
 - 3120 – $\mathbb{E}[\mathcal{E}_{i,j}] = 0$ for all $(i, j) \in [N] \times [d]$

3132 **Loss function**

- 3133
- 3134 • Define

$$3135 \quad \overline{L}_{ij}(\mathbf{u}, \boldsymbol{\theta}, \mathcal{E}) = (z_{ij} - y_{ij})^2 = (z_{ij} - f(\mathbf{u})_{ij} - \mathcal{E}_{ij})^2 \quad (49)$$

- 3137 • $L(\mathbf{u}) = \sum_{i,j} \overline{L}_{ij}(\mathbf{u}, \boldsymbol{\theta}, \mathcal{E})$
- 3139 • Training loss, $\overline{L}^{(t)}(\boldsymbol{\theta}) = \overline{L}^{(t)} = \mathbb{E}_{\substack{\mathbf{u} \sim \Delta \\ \mathcal{E} \sim \Delta_{\mathcal{E}}}} [L(\mathbf{u})]$
- 3141 • $\nabla_{\boldsymbol{\theta}} \overline{L}^{(t)}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{u}, \mathcal{E}} \sum_{i,j} \nabla_{\boldsymbol{\theta}} \overline{L}_{ij}(\mathbf{u})$

3142

3143 **The Goal** Given a target function f , what can we infer for $\boldsymbol{\theta} = (\mathbf{W}, \mathbf{K}, \mathbf{B}^{(1)}, \mathbf{B}^{(2)})$ from

3144 $\nabla_{\boldsymbol{\theta}} \overline{L}(\boldsymbol{\theta}) = \mathbf{0}$

- 3145
- 3146 1. Ideally, we'd like to assume that f can be represented exactly by 1-layer BASECONV.
 - 3147 2. For now, lets assume that $f(\mathbf{u})_{ij}$ only depends on \mathbf{u}_i :

3148 This includes as special cases:

- 3149 • $f(\mathbf{u})_{ij} = \mathbf{u} \odot \mathbf{u}$
- 3150 • $f(\mathbf{u}) = \mathbf{u} \cdot \overline{\mathbf{W}}$ for $\overline{\mathbf{W}} \in \mathbb{R}^{d \times d}$

3151

3152 We want to prove that there is a unique set of parameters that perform exactly these functions,

3153 resulting in the gradients of their loss function to be 0.

3154

3155 **D.6.2 A GENERIC PARTIAL DERIVATIVE**

3156 Lets try and reason as much as we can for a generic partial derivative. Let $x \in \boldsymbol{\theta} =$

3157 $(\mathbf{W}, \mathbf{K}, \mathbf{B}^{(1)}, \mathbf{B}^{(2)})$. Then from Equation (49), we have that for any i, j :

$$3158 \quad \frac{\partial \overline{L}_{i,j}}{\partial x} = 2(z_{ij} - f(\mathbf{u})_{ij} - \mathcal{E}_{ij}) \frac{\partial Z_{i,j}}{\partial x}$$

$$3159 \quad = 2 \left(\mathbf{T}_{ij}^{(1)} - \mathbf{T}_{ij}^{(2)} - \mathbf{T}_{ij}^{(3)} \right),$$

3160 where

$$3161 \quad \mathbf{T}_{ij}^{(1)} = Z_{i,j} \frac{\partial Z_{i,j}}{\partial x}.$$

$$3162 \quad \mathbf{T}_{ij}^{(2)} = f(u)_{ij} \frac{\partial Z_{i,j}}{\partial x}.$$

$$3163 \quad \mathbf{T}_{ij}^{(3)} = \mathcal{E}_{ij} \frac{\partial Z_{i,j}}{\partial x}.$$

3164 **Proposition D.47.** $\mathbb{E}_{\mathbf{u}, \mathcal{E}} [\mathbf{T}_{ij}^{(3)}] = 0$.

3165 *Proof.* Δ and $\Delta_{\mathcal{E}}$ are independent and $\mathbb{E}[\mathcal{E}_{ij}] = 0$ □

3166 From now on, we'll ignore the $\mathbf{T}_{ij}^{(3)}$ term because of Proposition D.47 we can (in expectation) assume

3167 that $\mathcal{E} = 0$.

3168

3169 **D.7 SETTING THE GRADIENTS TO 0**

3170 In this section we will prove the results of BASECONV ability to exactly compute SQUARE and

3171 LINEAR given assumptions on the input data and parameters. The results are as follows. First for the

3172 SQUARE function we have

3186 **Theorem D.48.** Given Assumptions [D.45](#), [D.46](#), [D.55](#), and a function

$$f(\mathbf{u}) = \mathbf{u} \odot \mathbf{u}.$$

3187
3188
3189 Let θ be such that, $\mathbb{E}\nabla_{\theta}\bar{L} = \mathbf{0}$ then $\text{BASECONV}(\mathbf{u}, \theta) = f(\mathbf{u})$.

3190
3191 And the following result for the LINEAR function,

3192 **Theorem D.49.** Given Assumptions [D.45](#), [D.46](#), [D.55](#), [D.61](#), and a function

$$f(\mathbf{u}) = \mathbf{u}\bar{W}.$$

3193
3194 Let θ be such that, $\mathbb{E}\nabla_{\theta}\bar{L} = \mathbf{0}$ then $\text{BASECONV}(\mathbf{u}, \theta) = f(\mathbf{u})$.

3195
3196 The following is to provide information about each entry in the output of a BASECONV layer.

3197
3198 **Lemma D.50.** For all $(i, j) \in [N] \times [d]$, the entries of a resulting layer of BASECONV, \mathbf{Z} , are:

$$3199 \quad \mathbf{Z}_{i,j} = \left(\left(\sum_{\ell=0}^{d-1} \mathbf{u}_{i,\ell} \cdot \mathbf{W}_{\ell,j} \right) + \mathbf{B}_{i,j}^{(1)} \right) \cdot \left(\left(\sum_{k=0}^i \mathbf{K}_{k,j} \cdot \mathbf{u}_{i-k,j} \right) + \mathbf{B}_{i,j}^{(2)} \right) \quad (50)$$

3200
3201 *Proof.* To begin, from [Equation \(48\)](#), we know that a layer of BASECONV yields a matrix \mathbf{Z} as,

$$3202 \quad \mathbf{Z} = \left(\mathbf{u} \cdot \mathbf{W} + \mathbf{B}^{(1)} \right) \odot \left(\left(\mathbf{K} * \mathbf{u} + \mathbf{B}^{(2)} \right) \right).$$

3203
3204 Looking at the $\mathbf{u} \cdot \mathbf{W}$ operation, we know that for a row $i \in [N]$ and column $j \in [d]$, the vector dot product is computed as

$$3205 \quad \langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle = \sum_{\ell=0}^{d-1} \mathbf{u}_{i,\ell} \cdot \mathbf{W}_{\ell,j}.$$

3206
3207 Meaning each entry in the resulting matrix is defined as such

$$3208 \quad (\mathbf{u} \cdot \mathbf{W})_{i,j} = \sum_{\ell=0}^{d-1} \mathbf{u}_{i,\ell} \cdot \mathbf{W}_{\ell,j}.$$

3209
3210 To sum the matrix $\mathbf{B}^{(1)}$ to this operation, we simply add the corresponding index giving us

$$3211 \quad \left(\mathbf{u} \cdot \mathbf{W} + \mathbf{B}^{(1)} \right)_{i,j} = \left(\sum_{\ell=0}^{d-1} \mathbf{u}_{i,\ell} \cdot \mathbf{W}_{\ell,j} \right) + \mathbf{B}_{i,j}^{(1)}. \quad (51)$$

3212
3213 Then, for the convolution operation between \mathbf{K} and \mathbf{u} , that's computed column by column, we have for all $j \in [d]$:

$$3214 \quad (\mathbf{K} * \mathbf{u})_{:,j} = \mathbf{K}_{:,j} * \mathbf{u}_{:,j},$$

3215
3216 i.e. for any $i \in [N]$,

$$3217 \quad (\mathbf{K}_{:,j} * \mathbf{u}_{:,j})[i] = \sum_{k=0}^i \mathbf{K}_{k,j} \cdot \mathbf{u}_{i-k,j}.$$

3218
3219 Finally, to sum $\mathbf{B}^{(2)}$ we add the corresponding entry giving us

$$3220 \quad \left(\mathbf{K} * \mathbf{u} + \mathbf{B}^{(2)} \right) = \left(\sum_{k=0}^i \mathbf{K}_{k,j} \cdot \mathbf{u}_{i-k,j} \right) + \mathbf{B}_{i,j}^{(2)}. \quad (52)$$

3221
3222 Combining [Equations \(51\)](#) and [\(52\)](#), gives us [Equation \(50\)](#) as expected. \square

3240 D.7.1 SOME PARTIAL DERIVATIVES ARE ALWAYS ZERO
3241

3242 To simplify future computations in this section we will state a simple lemma on chunks of the
3243 gradients of pieces of a layer that go to 0.

3244 **Lemma D.51.** Fix i, j . Then for any $j' \neq j$, $0 \leq \ell < N$, and $0 \leq k < N$ we have,
3245

$$3246 \frac{\partial \mathbf{Z}_{i,j'}}{\partial \mathbf{W}_{\ell,j}} = \frac{\partial \mathbf{Z}_{i,j'}}{\partial \mathbf{K}_{k,j}} = 0.$$

3247 Further, any $(i', j') \neq (i, j)$ we have,
3248

$$3249 \frac{\partial \mathbf{Z}_{i',j'}}{\partial \mathbf{B}_{i,j}^{(1)}} = \frac{\partial \mathbf{Z}_{i',j'}}{\partial \mathbf{B}_{i,j}^{(2)}} = 0.$$

3250 *Proof.* Follows from Equation (50) and definition of partial derivatives. \square
3251
3252

3253 D.7.2 GENERIC FORM OF PARTIAL DERIVATIVES PLUS A CONSEQUENCE
3254

3255 Given Lemma D.50 we can conclude the following.
3256

3257 **Lemma D.52.** For $0 \leq i < N$ and $0 \leq j < d$, any entry $x \in \{\mathbf{W}_{i,j}, \mathbf{B}_{i,j}^{(1)}\}$,
3258

$$3259 \frac{\partial \mathbf{Z}_{i,j}}{\partial x} = \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] + \mathbf{B}_{i,j}^{(2)} \right) \frac{\partial}{\partial x} \left(\langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle + \mathbf{B}_{i,j}^{(1)} \right)$$

3260 and $0 \leq i < N$ and $0 \leq j < d$, any entry $x \in \{\mathbf{K}_{i,j}, \mathbf{B}_{i,j}^{(2)}\}$,
3261

$$3262 \frac{\partial \mathbf{Z}_{i,j}}{\partial x} = \left(\langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle + \mathbf{B}_{i,j}^{(1)} \right) \cdot \frac{\partial}{\partial x} \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] + \mathbf{B}_{i,j}^{(2)} \right).$$

3263 A consequence of Lemma D.52 is the following.
3264

3265 **Corollary D.53.** Let $\theta = (\mathbf{W}, \mathbf{K}, \mathbf{B}^{(1)}, \mathbf{B}^{(2)}) = \mathbf{0}$. Then for all parameter variable's x , we have
3266

$$3267 \frac{\partial \mathbf{Z}_{i,j}}{\partial x} = 0.$$

3268 Specifically,
3269

$$3270 \nabla_{\theta} \bar{L}(\theta) |_{\theta=\mathbf{0}} = \mathbf{0}.$$

3271 Corollary D.53 implies that initializing $\theta = \mathbf{0}$ is not a good choice for initializing parameters since it
3272 is a local minima.
3273

3274 We can exactly figure out the partial derivatives of Lemma D.52 by the following.
3275

3276 **Lemma D.54.** Fix i, j . Then for any $0 \leq \ell < N$ we have,
3277

$$3278 \frac{\partial}{\partial \mathbf{W}_{\ell,j}} \left(\langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle + \mathbf{B}_{i,j}^{(1)} \right) = \mathbf{u}_{i,\ell}$$

3279 and
3280

$$3281 \frac{\partial}{\partial \mathbf{B}_{i,j}^{(1)}} \left(\langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle + \mathbf{B}_{i,j}^{(1)} \right) = 1.$$

3282 Also,
3283

$$3284 \frac{\partial}{\partial \mathbf{B}_{i,j}^{(2)}} \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] + \mathbf{B}_{i,j}^{(2)} \right) = 1.$$

3285 Next, for any $0 \leq k \leq i$, we have
3286

$$3287 \frac{\partial}{\partial \mathbf{K}_{k,j}} \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] + \mathbf{B}_{i,j}^{(2)} \right) = \mathbf{u}_{i-k,j}$$

3294 and for all $k > i$,

$$3295 \frac{\partial}{\partial \mathbf{K}_{k,j}} \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] + \mathbf{B}_{i,j}^{(2)} \right) = 0.$$

3297 We don't need to worry about $k > i$ since those pieces will also be 0 due to the definition of
3298 convolution.

3300 *Proof.* Let's begin by looking at

$$3301 \frac{\partial}{\partial \mathbf{W}_{\ell,j}} \left(\langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle + \mathbf{B}_{i,j}^{(1)} \right).$$

3302 Expanding this out gives us

$$3303 \frac{\partial}{\partial \mathbf{W}_{\ell,j}} \left(\sum_{\ell=0}^{d-1} \mathbf{u}_{i,\ell} \cdot \mathbf{W}_{\ell,j} + \mathbf{B}_{i,j}^{(1)} \right).$$

3308 When we take the partial derivative of this with respect to $\mathbf{W}_{\ell,j}$, $\mathbf{B}_{i,j}^{(1)}$ goes to 0. And the term

$$3309 \frac{\partial}{\partial \mathbf{W}_{\ell,j}} \left(\sum_{\ell=0}^{d-1} \mathbf{u}_{i,\ell} \cdot \mathbf{W}_{\ell,j} \right) = \mathbf{u}_{i,\ell},$$

3312 as desired.

3313 Next, let us look at

$$3314 \frac{\partial}{\partial \mathbf{B}_{i,j}^{(1)}} \left(\langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle + \mathbf{B}_{i,j}^{(1)} \right)$$

3318 Since $\mathbf{B}_{i,j}^{(1)}$ doesn't show up in the dot product of the vectors, we know that piece goes to zero, giving
3319 us

$$3320 \frac{\partial}{\partial \mathbf{B}_{i,j}^{(1)}} \left(\langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle + \mathbf{B}_{i,j}^{(1)} \right) = \frac{\partial \mathbf{B}_{i,j}^{(1)}}{\partial \mathbf{B}_{i,j}^{(1)}} = 1,$$

3323 as desired.

3324 Next, for any $0 \leq k \leq i$ we have

$$3325 \frac{\partial}{\partial \mathbf{K}_{k,j}} \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] + \mathbf{B}_{i,j}^{(2)} \right).$$

3328 The $\mathbf{B}_{i,j}^{(2)}$ term goes to 0 as we're taking the partial derivative with respect to $\mathbf{K}_{k,j}$. So we have

$$3329 \frac{\partial}{\partial \mathbf{K}_{k,j}} \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] \right) = \frac{\partial}{\partial \mathbf{K}_{k,j}} \sum_{k'=0}^i \mathbf{K}_{k',j} \mathbf{u}_{i-k',j} = \mathbf{u}_{i-k,j}$$

3332 as desired, since $\mathbf{K}_{k,j}$ only shows up in the summation when $k' = k$. Then for $k > i$ we have

$$3333 \frac{\partial}{\partial \mathbf{K}_{k,j}} \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] \right) = \frac{\partial}{\partial \mathbf{K}_{k,j}} \sum_{k'=0}^i \mathbf{K}_{k',j} \mathbf{u}_{i-k',j} = 0 \quad (53)$$

3336 as desired, since $\mathbf{K}_{k,j}$ will never show up in the summation as $k' < k$.

3337 Finally, let's look at the fourth piece,

$$3338 \frac{\partial}{\partial \mathbf{B}_{i,j}^{(2)}} \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] + \mathbf{B}_{i,j}^{(2)} \right).$$

3341 The term $\mathbf{B}_{i,j}^{(2)}$ doesn't appear in the result of the convolution operation, therefore that piece goes to
3342 0, giving us

$$3343 \frac{\partial}{\partial \mathbf{B}_{i,j}^{(2)}} \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] + \mathbf{B}_{i,j}^{(2)} \right) = \frac{\partial \mathbf{B}_{i,j}^{(2)}}{\partial \mathbf{B}_{i,j}^{(2)}} = 1,$$

3346 as desired.

3347

□

Another assumption Before our main result we have another restriction on the target f which is that f must be defined with a linear map, $\overline{\mathbf{W}} \in \mathbb{R}^{d \times d}$ that has non-zero columns. We make further assumptions on \mathbf{W} , $\mathbf{B}^{(1)}$, \mathbf{K} , $\mathbf{B}^{(2)}$ which we will justify later. We note that this assumption is satisfied for both the SQUARE and LINEAR functions.

Assumption D.55. For all j either $(\mathbf{K}_{:,j} \neq \mathbf{0})$ or $(\mathbf{B}_{:,j}^{(2)} \neq \mathbf{0})$ and $\mathbf{W}_{:,j} \neq \mathbf{0}$. Further, we have $\mathbf{B}^{(1)} = \mathbf{0}$.

The target function $f : \mathbb{R}^{N \times d} \rightarrow \mathbb{R}^{N \times d}$ is

1. 1-layer BC $(\overline{\mathbf{W}}, \overline{\mathbf{K}}, \overline{\mathbf{B}}^{(1)}, \overline{\mathbf{B}}^{(2)})$ such that for all j , $\|\overline{\mathbf{W}}_{:,j}\|_2 \geq 0$
2. $f(\mathbf{u})_{i,j}$ only depends on $\mathbf{u}_{i,:}$

With the above assumptions and definitions, we are finally ready to look at the expected partial derivatives of pieces of the loss function.

Definition D.56. For the rest of the section, we will redefine $\theta = (\mathbf{W}, \mathbf{K}, \mathbf{B}^{(2)})$. Note that we are just removing $\mathbf{B}^{(1)}$ since it is all zeros as per [Assumption D.55](#).

Lemma D.57. Given [Assumption D.45](#) and recall that $\mathbf{B}^{(1)} = \mathbf{0}$. Fix i, j . Then we have

$$\mathbb{E} \left[\mathbf{Z}_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] = \mathbf{B}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'}^2] \mathbf{W}_{\ell',j}^2.$$

Next, for any $0 \leq k \leq i$ we have

$$\mathbb{E} \left[\mathbf{Z}_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = \mathbf{K}_{k,j} \sum_{\ell'=0}^{d-1} \mathbf{W}_{\ell',j}^2 \mathbb{E} [\mathbf{u}_{i,\ell'}^2] \mathbb{E} [\mathbf{u}_{i-k,j}^2].$$

(Recall for $k > i$ partial derivatives are 0).

Finally, for any $0 \leq \ell \leq d-1$,

$$\mathbb{E} \left[\mathbf{Z}_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] = \mathbf{W}_{\ell,j} \sum_{k'=0}^i \mathbf{K}_{k',j}^2 \mathbb{E} [\mathbf{u}_{i-k',j}^2 \mathbf{u}_{i,\ell}^2] + (\mathbf{B}_{i,j}^{(2)})^2 \mathbf{W}_{\ell,j} \mathbb{E} [\mathbf{u}_{i,\ell}^2].$$

Proof. Given [Lemma D.52](#) and [Lemma D.54](#) (along with the fact that $\mathbf{B}^{(1)} = \mathbf{0}$) we have

$$\begin{aligned} \mathbb{E} \left[\mathbf{Z}_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] &= \mathbb{E} \left[\left(\mathbf{K}_{:,j} * \mathbf{u}_{:,j}[i] + \mathbf{B}_{i,j}^{(2)} \right) \langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle^2 \right] \\ &= \sum_{\ell'=0}^{d-1} \sum_{\ell''=0}^{d-1} \sum_{k'=0}^i \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i-k',j} \mathbf{u}_{i,\ell''}] \mathbf{W}_{\ell',j} \mathbf{W}_{\ell'',j} \mathbf{K}_{k',j} \\ &\quad + \mathbf{B}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \sum_{\ell''=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell''}] \mathbf{W}_{\ell',j} \mathbf{W}_{\ell'',j}. \end{aligned}$$

In the above, the first summation goes to 0 since for all ℓ', ℓ'', k , by [Assumption D.45](#), as the expected value of the product of three \mathbf{u} 's will be 0 since there's an odd number of them. Again, by [Assumption D.45](#), the second summation will be non-zero if and only if $\ell' = \ell''$. Therefore we get the following,

$$\mathbb{E} \left[\mathbf{Z}_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] = \mathbf{B}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'}^2] \mathbf{W}_{\ell',j}^2$$

as desired.

Moving onto the next piece, using [Lemma D.52](#) and [Lemma D.54](#) (along with the fact that $B^{(1)} = \mathbf{0}$) we have

$$\begin{aligned} \mathbb{E} \left[\mathbf{Z}_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] &= \mathbb{E} \left[\left(\mathbf{K}_{:,j} * \mathbf{u}_{:,j}[i] + \mathbf{B}_{i,j}^{(2)} \right) \langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle^2 \mathbf{u}_{i-k,j} \right] \\ &= \sum_{\ell'=0}^{d-1} \sum_{\ell''=0}^{d-1} \sum_{k'=0}^i \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell''} \mathbf{u}_{i-k',j} \mathbf{u}_{i-k,j}] \mathbf{W}_{\ell',j} \mathbf{W}_{\ell'',j} \mathbf{K}_{k',j} \\ &\quad + \mathbf{B}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \sum_{\ell''=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell''} \mathbf{u}_{i-k,j}] \mathbf{W}_{\ell',j} \mathbf{W}_{\ell'',j} \end{aligned}$$

By [Assumption D.45](#), only expected values of terms with square monomials are non-zero. Specifically, the first summation has the $\mathbf{u}_{i-k,j}$ term, therefore, we need $k' = k$ to get an even exponent. This is the same reasoning for $\ell' = \ell''$. Therefore, the first summation is non-zero if and only if $k' = k$ and $\ell' = \ell''$. The second summation will be 0 since for all ℓ', ℓ'' the expected value of the \mathbf{u} 's is 0 since there's an odd number of them, there will always be an odd exponent. So we get

$$\mathbb{E} \left[\mathbf{Z}_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = \mathbf{K}_{k,j} \sum_{\ell'=0}^{d-1} \mathbf{W}_{\ell',j}^2 \mathbb{E} [\mathbf{u}_{i,\ell'}^2] \mathbb{E} [\mathbf{u}_{i-k,j}^2]$$

as desired. In the above we note that by [Assumption D.45](#), any two entries in \mathbf{u} are independent random variables so we can multiply their expectations as above.

Moving onto the final piece, given [Lemma D.52](#) and [Lemma D.54](#) and $B^{(1)} = \mathbf{0}$ we have

$$\begin{aligned} \mathbb{E} \left[\mathbf{Z}_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] &= \mathbb{E} \left[\left(\mathbf{K}_{:,j} * \mathbf{u}_{:,j} \right) [i] + \left(\mathbf{B}_{i,j}^{(2)} \right) \right]^2 \langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle \mathbf{u}_{i,\ell} \\ &= \sum_{k'=0}^i \sum_{k''=0}^i \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i-k',j} \mathbf{u}_{i,\ell'} \mathbf{u}_{i-k'',j} \mathbf{u}_{i,\ell}] \mathbf{K}_{k',j} \mathbf{K}_{k'',j} \mathbf{W}_{\ell',j} \\ &\quad + 2\mathbf{B}_{i,j}^{(2)} \sum_{k'=0}^i \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i-k',j} \mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell}] \mathbf{K}_{k',j} \mathbf{W}_{\ell',j} \\ &\quad + \left(\mathbf{B}_{i,j}^{(2)} \right)^2 \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell}] \mathbf{W}_{\ell',j}. \end{aligned}$$

We again use [Assumption D.45](#) to simplify the summations. The first summation has the $\mathbf{u}_{i,\ell}$ term, therefore to get an even exponent on it we need $\ell' = \ell$. This is the same reasoning for $k' = k''$. Therefore the first summation will be non-zero if and only if $\ell' = \ell$ and $k' = k''$. The second summation will be 0 for all k', ℓ' since we're taking the expected value of an odd number of \mathbf{u} products, there will always be an odd exponent. The third term will be non-zero if and only if $\ell' = \ell$ to get an even exponent on \mathbf{u} 's entry. Therefore we have,

$$\mathbb{E} \left[\mathbf{Z}_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] = \mathbf{W}_{\ell,j} \sum_{k'=0}^i \mathbf{K}_{k',j}^2 \mathbb{E} [\mathbf{u}_{i-k',j}^2] \mathbb{E} [\mathbf{u}_{i,\ell}^2] + \left(\mathbf{B}_{i,j}^{(2)} \right)^2 \mathbf{W}_{\ell,j} \mathbb{E} [\mathbf{u}_{i,\ell}^2]$$

as desired. In the above we note that by [Assumption D.45](#), any two entries in \mathbf{u} are independent random variables. □

Lemma D.58. Fix i, j . Then we have

$$\mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}_{i,j}^{(2)}} \right] = 0.$$

Next, we have

$$\mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{0,j}} \right] = \mathbf{W}_{j,j} \mathbb{E} [\mathbf{u}_{i,j}^4].$$

3456 For $k > 0$, we have

$$3457 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = 0.$$

3459 Next,

$$3460 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{j,j}} \right] = \mathbf{K}_{0,j} \mathbb{E} [\mathbf{u}_{i,j}^4].$$

3462 otherwise, when $\ell \neq j$ we have

$$3463 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] = 0.$$

3466 *Proof.* Given [Lemma D.52](#) and [Lemma D.54](#) we have

$$3467 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] = \langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle \mathbf{u}_{i,j}^2$$

$$3470 = \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,j}^2] \mathbf{W}_{\ell',j}.$$

3473 We simplify the above using [Assumption D.45](#). Therefore, the summation is 0 for all ℓ since we're taking the expected value of an odd number of \mathbf{u} 's. Therefore,

$$3475 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] = 0.$$

3478 Next, by [Lemma D.52](#) and [Lemma D.54](#) we have

$$3479 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = \langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle \mathbf{u}_{i-k,j} \mathbf{u}_{i,j}^2$$

$$3482 = \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,j}^2 \mathbf{u}_{i-k,j} \mathbf{u}_{i,\ell'}] \mathbf{W}_{\ell',j}.$$

3485 We simplify the above using [Assumption D.45](#). The summation will be non-zero only when $k = 0$ and $\ell' = j$ to ensure we have an even exponent on the \mathbf{u} variable, since that is the only way to get even exponents on the \mathbf{u} 's. Therefore we have

$$3488 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{0,j}} \right] = \mathbf{W}_{j,j} \mathbb{E} [\mathbf{u}_{i,j}^4],$$

3490 and for $k > 0$,

$$3491 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = 0.$$

3494 Moving onto the final piece, given [Lemma D.52](#) and [Lemma D.54](#), we have

$$3495 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] = \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}[i] + \mathbf{B}_{i,j}^{(2)}) \right) \mathbf{u}_{i,\ell} \mathbf{u}_{i,j}^2$$

$$3498 = \sum_{k'=0}^i \mathbb{E} [\mathbf{u}_{i-k',j} \mathbf{u}_{i,\ell} \mathbf{u}_{i,j}^2] \mathbf{K}_{k',j} + \mathbf{B}_{i,j}^{(2)} \mathbb{E} [\mathbf{u}_{i,\ell} \mathbf{u}_{i,j}^2].$$

3501 We simplify the above using [Assumption D.45](#). The summation will be non-zero only when $k' = 0$ and $\ell = j$, since that is the only way to get even exponents on the \mathbf{u} 's. The second term will be 0 since we're taking the expected value of a non-square monomial. Therefore we have,

$$3504 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{j,j}} \right] = \mathbf{K}_{0,j} \mathbb{E} [\mathbf{u}_{i,j}^4].$$

3506 Otherwise, when $\ell \neq j$ we have,

$$3507 \mathbb{E} \left[\mathbf{u}_{i,j}^2 \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] = 0$$

3509 as desired. □

Theorem D.59. Given Assumptions D.45, D.46, D.55, and a function

$$f(\mathbf{u}) = \mathbf{u} \odot \mathbf{u}.$$

Let $\boldsymbol{\theta}$ be such that, $\mathbb{E}\nabla_{\boldsymbol{\theta}}\bar{L} = \mathbf{0}$ then $\text{BASECONV}(\mathbf{u}, \boldsymbol{\theta}) = f(\mathbf{u})$.

Proof. Recall our loss function from Equation (49). Then given Lemma D.57 and Lemma D.58 we know that for any $k \geq 0$,

$$\begin{aligned} \mathbb{E} \left[\frac{\partial \bar{L}}{\partial \mathbf{K}_{k,j}} \right] &= \sum_{i',j'} \mathbb{E} \left[\frac{\partial \bar{L}_{i',j'}}{\partial \mathbf{K}_{k,j}} \right] \\ &= \sum_{i'} \mathbb{E} \left[(\mathbf{Z}_{i',j} - \mathbf{u}_{i',j}^2) \frac{\partial \mathbf{Z}_{i',j}}{\partial \mathbf{K}_{k,j}} \right]. \end{aligned}$$

In the above, the second equality follows from Lemma D.51.

From Lemma D.57, Lemma D.58 we have the following for $k > 0$,

$$\mathbb{E} \left[\frac{\partial \bar{L}}{\partial \mathbf{K}_{k,j}} \right] = \mathbf{K}_{k,j} \sum_{i'} \sum_{\ell'=0}^{d-1} \mathbf{W}_{\ell',j}^2 \mathbb{E} [\mathbf{u}_{i',\ell'}^2] \mathbb{E} [\mathbf{u}_{i'-k,j}^2].$$

We know that from Assumption D.46 that the expected value of $\mathbb{E} [\mathbf{u}_{i',\ell'}^2]$ and $\mathbb{E} [\mathbf{u}_{i'-k,j}^2]$ are strictly positive due to having an even exponent. Also, due to Assumption D.55 we know that $\mathbf{W}_{:,j}$ has at least one non-zero entry for all j . This implies that for at least one ℓ' , $\mathbf{W}_{\ell',j}^2 > 0$ (and hence the summation is strictly a positive value). Therefore to set this partial derivative to 0 implies that $\mathbf{K}_{k,j} = 0$. Explicitly, $\mathbf{K}_{k,j} = 0$ for all j and $k > 0$. This in turn implies that for all j :

$$\mathbf{K}_{:,j} \neq \mathbf{0} \Leftrightarrow \mathbf{K}_{0,j} \neq 0. \quad (54)$$

Next we'll examine the loss function when we take the partial derivative with respect to $\mathbf{W}_{\ell,j}$. For all $\ell \neq j$,

$$\begin{aligned} \mathbb{E} \left[\frac{\partial \bar{L}}{\partial \mathbf{W}_{\ell,j}} \right] &= \sum_{i',j'} \mathbb{E} \left[\frac{\partial \bar{L}_{i',j'}}{\partial \mathbf{W}_{\ell,j}} \right] \\ &= \sum_{i'} \left(\mathbf{W}_{\ell,j} \sum_{k'=0}^{i'} \mathbf{K}_{k',j}^2 \mathbb{E} [\mathbf{u}_{i'-k',j}^2 \mathbf{u}_{i',\ell}^2] + \left(\mathbf{B}_{i',j}^{(2)} \right)^2 \mathbf{W}_{\ell,j} \mathbb{E} [\mathbf{u}_{i',\ell}^2] \right). \end{aligned}$$

In the above the second equality follows from Lemma D.51, Lemma D.57, and Lemma D.58.

We know from above that for $k > 0$, $\mathbf{K}_{k,j} = 0$. Thus, we can simplify the above to

$$\begin{aligned} &\sum_{i'} \mathbf{W}_{\ell,j} \left(\mathbf{K}_{0,j}^2 \mathbb{E} [\mathbf{u}_{i',j}^2] \mathbb{E} [\mathbf{u}_{i',\ell}^2] + \left(\mathbf{B}_{i',j}^{(2)} \right)^2 \mathbb{E} [\mathbf{u}_{i',\ell}^2] \right) \\ &= \mathbf{W}_{\ell,j} \sum_{i'} \mathbb{E} [\mathbf{u}_{i',\ell}^2] \left(\mathbf{K}_{0,j}^2 \mathbb{E} [\mathbf{u}_{i',j}^2] + \left(\mathbf{B}_{i',j}^{(2)} \right)^2 \right). \end{aligned}$$

By Equation (54) we know $\mathbf{K}_{0,j} > 0$ and from Assumption D.55 at least one entry in each column of $\mathbf{B}^{(2)}$ may be non-zero, $\mathbf{B}_{:,j}^{(2)} > 0$. Therefore

$$\left(\mathbf{K}_{0,j}^2 \mathbb{E} [\mathbf{u}_{i',j}^2] + \mathbf{B}_{i',j}^{(2)} \right) \neq 0 \implies \mathbf{W}_{\ell,j} = 0$$

Since we proved that $\mathbf{W}_{\ell,j} = 0$ for all $\ell \neq j$, and we know at least one entry in each column $\mathbf{W}_{:,j} \neq \mathbf{0}$, we must have

$$\mathbf{W}_{j,j} \neq 0. \quad (55)$$

3564 Moving onto $\mathbf{B}_{i,j}^{(2)}$, we have

$$3566 \mathbb{E} \left[\frac{\partial \bar{L}}{\partial \mathbf{B}_{i,j}^{(2)}} \right] = \mathbf{B}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \mathbf{W}_{\ell',j}^2 \mathbb{E} [\mathbf{u}_{i,\ell'}^2].$$

3569 The above summation piece is not equal to zero as we know $\mathbf{W}_{j,j} \neq 0$ and hence, $\mathbf{W}_{j,j}^2 > 0$. Also
 3570 note the exponents on the \mathbf{u} 's is even, and so by [Assumption D.45](#) we have for all ℓ' , $\mathbb{E} [\mathbf{u}_{i,\ell'}^2] > 0$.
 3571 Therefore, if we set

$$3573 \mathbb{E} \left[\frac{\partial \bar{L}}{\partial \mathbf{B}_{i,j}^{(2)}} \right] = 0$$

3574 we get for all i, j ,

$$3577 \mathbf{B}_{i,j}^{(2)} = 0.$$

3578 Now that we know $\mathbf{B}^{(2)} = \mathbf{0}$ this will help us simplify the following. Given [Lemma D.57](#) and
 3579 [Lemma D.58](#), for $\ell = j$, we have

$$3581 \mathbb{E} \left[\frac{\partial \bar{L}}{\partial \mathbf{W}_{j,j}} \right] = \sum_{i'} \mathbf{W}_{j,j} \mathbf{K}_{0,j}^2 \mathbb{E} [\mathbf{u}_{i',j}^4] - \mathbf{K}_{0,j} \mathbb{E} [\mathbf{u}_{i',j}^4],$$

3584 we simplify and remove $\mathbf{B}^{(2)}$ terms and get

$$3586 \mathbf{K}_{0,j} (\mathbf{K}_{0,j} \mathbf{W}_{j,j} - 1) \left(\sum_{i'} \mathbb{E} [\mathbf{u}_{i',j}^4] \right).$$

3589 Setting this to 0 we can conclude that since $\mathbf{K}_{0,j} \neq 0$ and the above summation isn't 0, we get
 3590 $\mathbf{K}_{0,j} \mathbf{W}_{j,j} = 1$.

3591 So far we have that $\mathbf{B}^{(2)} = \mathbf{0}$ and $\mathbf{K}_{k,j} = 0$ for $k > 0$. Also that $\mathbf{K}_{0,j} \neq 0$ for all j . Then by
 3592 [Lemma D.58](#) and [Lemma D.57](#), consider

$$3594 \mathbb{E} \left[\frac{\partial \bar{L}}{\partial \mathbf{K}_{0,j}} \right] = \sum_{i'} (\mathbf{K}_{0,j} \mathbf{W}_{j,j}^2 \mathbb{E} [\mathbf{u}_{i',j}^2 \mathbf{u}_{i',j}^2] - \mathbf{W}_{j,j} \mathbb{E} [\mathbf{u}_{i',j}^4])$$

$$3597 = \mathbf{W}_{j,j} (\mathbf{K}_{0,j} \mathbf{W}_{j,j} - 1) \left(\sum_{i'} \mathbb{E} [\mathbf{u}_{i',j}^4] \right). \quad (56)$$

3599 We know from (55) that $\mathbf{W}_{j,j} \neq 0$ and the summation piece isn't 0 by [Assumption D.45](#). Therefore,

$$3601 \mathbb{E} \left[\frac{\partial \bar{L}}{\partial \mathbf{K}_{0,j}} \right] = 0 \implies \mathbf{K}_{0,j} \cdot \mathbf{W}_{j,j} - 1 = 0$$

$$3603 \implies \mathbf{K}_{0,j} \cdot \mathbf{W}_{j,j} = 1$$

3604 as desired.

3606 Note that this implies for all j :

$$3608 \mathbf{K}_{0,j} = \frac{1}{\mathbf{W}_{j,j}}.$$

3610 Given [Assumption D.55](#) and the above values; we get that \mathbf{W} is a diagonal matrix, $\mathbf{K} =$
 3611 $\begin{pmatrix} \mathbf{w}_{0,0} & \mathbf{w}_{1,1} & \dots & \mathbf{w}_{d-1,d-1} \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \end{pmatrix}$, $\mathbf{B}^{(2)} = \mathbf{0}^{N \times d}$. Lets use these pieces to show that $\text{BASECONV}(\mathbf{u}) = \mathbf{u} \odot \mathbf{u}$
 3612 (recall that $\mathbf{B}^{(1)} = \mathbf{0}$).

3614
 3615 Indeed,

$$3617 \text{BASECONV}(\mathbf{u}) = (\mathbf{u} \cdot \mathbf{W}) \odot \left(\mathbf{K} * \mathbf{u} + \mathbf{B}^{(2)} \right)$$

3618 Plugging in our values we get

$$3619 \left(\mathbf{u} \cdot \begin{pmatrix} \mathbf{W}_{0,0} & 0 & \dots & 0 \\ 0 & \mathbf{W}_{1,1} & \dots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \dots & \mathbf{W}_{d-1,d-1} \end{pmatrix} \right) \odot \left(\left(\begin{pmatrix} \frac{1}{\mathbf{W}_{0,0}} & \frac{1}{\mathbf{W}_{1,1}} & \dots & \frac{1}{\mathbf{W}_{d-1,d-1}} \\ \mathbf{0}_{N-1 \times d-1} \end{pmatrix} * \mathbf{u} + \mathbf{0}^{N \times d} \right) \right).$$

3624 Note that

$$3625 \mathbf{K} * \mathbf{u} = \mathbf{u} \begin{pmatrix} \frac{1}{\mathbf{W}_{0,0}} & 0 & \dots & 0 \\ 0 & \frac{1}{\mathbf{W}_{1,1}} & \dots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \dots & \frac{1}{\mathbf{W}_{d-1,d-1}} \end{pmatrix} = \mathbf{u} \mathbf{W}^{-1},$$

3631 which gives us

$$3632 (\mathbf{u} \mathbf{W} \odot \mathbf{K} * \mathbf{u}) = (\mathbf{u} \mathbf{W} \odot \mathbf{u} \mathbf{W}^{-1}) = \mathbf{u} \odot \mathbf{u},$$

3633 where the last equality follows since \mathbf{W} is a diagonal matrix with all non-zero entries in it's diagonal. \square

3637 The following lemma will be for when the function we are considering is an linear map.

3638 **Lemma D.60.** Fix $\overline{\mathbf{W}} \in \mathbb{R}^{d \times d}$ and i, j . Then we have

$$3639 \mathbb{E} \left[\langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] = \sum_{\ell}^{d-1} \mathbf{W}_{\ell,j} \overline{\mathbf{W}}_{\ell,j} \mathbb{E} [\mathbf{u}_{i,\ell}^2].$$

3643 For all k ,

$$3644 \mathbb{E} \left[\langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = 0.$$

3646 For all ℓ ,

$$3647 \mathbb{E} \left[\langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] = \mathbf{B}_{i,j}^{(2)} \overline{\mathbf{W}}_{\ell,j} \mathbb{E} [\mathbf{u}_{i,\ell}^2].$$

3651 *Proof.* Given [Lemma D.52](#) and [Lemma D.54](#) (and the fact that $\mathbf{B}^{(1)} = \mathbf{0}$) we have

$$3652 \mathbb{E} \left[\langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] = \mathbb{E} [\langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle \langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle] \\ 3653 = \sum_{\ell=0}^{d-1} \sum_{\ell'=0}^{d-1} \mathbf{W}_{\ell,j} \overline{\mathbf{W}}_{\ell',j} \mathbb{E} [\mathbf{u}_{i,\ell} \mathbf{u}_{i,\ell'}]$$

3654 From [Assumption D.45](#) we get that the summation will always be non-zero if and only if $\ell = \ell'$ so that the \mathbf{u} variable has an even exponent. Therefore we get,

$$3655 \mathbb{E} \left[\langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] = \sum_{\ell=0}^{d-1} \mathbf{W}_{\ell,j} \overline{\mathbf{W}}_{\ell,j} \mathbb{E} [\mathbf{u}_{i,\ell}^2]$$

3661 as desired.

3662 Next, for all k , by [Lemma D.52](#) and [Lemma D.54](#) (and the fact that $\mathbf{B}^{(1)} = \mathbf{0}$)

$$3663 \mathbb{E} \left[\langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = \mathbb{E} [\langle \mathbf{u}_{i,:}, \mathbf{W}_{:,j} \rangle \mathbf{u}_{i-k,j} \langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle] \\ 3664 = \sum_{\ell=0}^{d-1} \sum_{\ell'=0}^{d-1} \mathbf{W}_{\ell,j} \overline{\mathbf{W}}_{\ell',j} \mathbb{E} [\mathbf{u}_{i,\ell} \mathbf{u}_{i,\ell'} \mathbf{u}_{i-k,j}].$$

We simplify the above using [Assumption D.45](#). The summation goes to 0 since there are an odd number of \mathbf{u} terms, there will always be an odd exponent. Note that this is true for $k \leq i$. Recall from [Equation \(53\)](#) that for $k > i$,

$$\frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} = 0.$$

Therefore, for all k ,

$$\mathbb{E} \left[\langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = 0$$

as desired.

Next, for all ℓ , by [Lemma D.52](#) and [Lemma D.54](#) (and the fact that $\mathbf{B}^{(1)} = \mathbf{0}$)

$$\begin{aligned} \mathbb{E} \left[\langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] &= \mathbb{E} \left[\left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] + \mathbf{B}_{i,j}^{(2)} \right) \mathbf{u}_{i,\ell} \langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle \right] \\ &= \left(\sum_{k'=0}^i \sum_{\ell'=0}^{d-1} \mathbf{K}_{k',j} \overline{\mathbf{W}}_{\ell',j} \mathbb{E} [\mathbf{u}_{i-k',j} \mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell}] \right) \\ &\quad + \mathbf{B}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \overline{\mathbf{W}}_{\ell',j} \mathbb{E} [\mathbf{u}_{i,\ell} \mathbf{u}_{i,\ell'}] \end{aligned}$$

We simplify the above using [Assumption D.45](#). The first summation will always be 0 due to an odd exponent on the \mathbf{u} 's. The second summation piece will always be non-zero if and only if $\ell' = \ell$, giving us the even exponent on the \mathbf{u} variable. Therefore we get,

$$\mathbb{E} \left[\langle \mathbf{u}_{i,:}, \overline{\mathbf{W}}_{:,j} \rangle \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] = \mathbf{B}_{i,j}^{(2)} \overline{\mathbf{W}}_{\ell,j} \mathbb{E} [\mathbf{u}_{i,\ell}^2]$$

as desired. \square

We make another assumption to assist with the following theorem on linear maps.

Assumption D.61. For all j , $\langle \mathbf{W}_{:,j}, \overline{\mathbf{W}}_{:,j} \rangle \neq 0$ and $\overline{\mathbf{W}}_{:,j} \neq \mathbf{0}$.

Theorem D.62. Given Assumptions [D.45](#), [D.46](#), [D.55](#), [D.61](#), and a function

$$f(\mathbf{u}) = \mathbf{u} \overline{\mathbf{W}}.$$

Let $\boldsymbol{\theta}$ be such that, $\mathbb{E} \nabla_{\boldsymbol{\theta}} \overline{\mathbf{L}} = \mathbf{0}$ then $\text{BASECONV}(\mathbf{u}, \boldsymbol{\theta}) = f(\mathbf{u})$.

Proof. From [lemma D.51](#) we get that

$$\begin{aligned} \mathbb{E} \left[\frac{\partial \overline{\mathbf{L}}}{\partial \mathbf{B}_{i,j}^{(2)}} \right] &= \sum_{i',j'} \mathbb{E} \left[\frac{\partial \overline{\mathbf{L}}_{i',j'}}{\partial \mathbf{B}_{i',j'}^{(2)}} \right] \\ &= \sum_{i,j} \mathbb{E} \left[\frac{\partial \overline{\mathbf{L}}_{i,j}}{\partial \mathbf{B}_{i,j}^{(2)}} \right] \end{aligned}$$

Recall our loss function from [Equation \(49\)](#). Then given [Lemma D.57](#) and [Lemma D.60](#) we know that

$$\sum_{i,j} \mathbb{E} \left[\frac{\partial \overline{\mathbf{L}}_{i,j}}{\partial \mathbf{B}_{i,j}^{(2)}} \right] = \mathbf{B}_{i,j}^{(2)} \sum_{\ell=0}^{d-1} \mathbf{W}_{\ell,j}^2 \mathbb{E} [\mathbf{u}_{i,\ell}^2] - \sum_{\ell=0}^{d-1} \mathbf{W}_{\ell,j} \overline{\mathbf{W}}_{\ell,j} \mathbb{E} [\mathbf{u}_{i,j}^2] = 0,$$

Which implies that

$$\mathbf{B}_{i,j}^{(2)} = \frac{\sum_{\ell=0}^{d-1} \mathbf{W}_{\ell,j} \overline{\mathbf{W}}_{\ell,j}}{\sum_{\ell=0}^{d-1} \mathbf{W}_{\ell,j}^2}.$$

Given [Assumption D.61](#) we know that the numerator will be non-zero and given assumption [Assumption D.55](#) we know the denominator will always be non-zero as well. Therefore we get that for all i, j

$$\mathbf{B}_{i,j} = \frac{\langle \mathbf{W}_{:,j}, \overline{\mathbf{W}}_{:,j} \rangle}{\langle \mathbf{W}_{:,j}, \mathbf{W}_{:,j} \rangle} \stackrel{\text{def}}{=} b_j. \quad (57)$$

Next, via [Lemma D.51](#) we have for all $k \geq 0$:

$$\begin{aligned} \mathbb{E} \left[\frac{\partial \overline{\mathcal{L}}}{\partial \mathbf{K}_{k,j}} \right] &= \sum_{i',j'} \mathbb{E} \left[\frac{\partial \overline{\mathcal{L}}_{i',j'}}{\partial \mathbf{K}_{k,j'}} \right] \\ &= \sum_{i'} \mathbb{E} \left[\frac{\partial \overline{\mathcal{L}}_{i',j}}{\partial \mathbf{K}_{k,j}} \right] \end{aligned}$$

Then from [Lemma D.57](#) and [lemma D.60](#) we get

$$\sum_{i'} \mathbb{E} \left[\mathbf{z}_{i',j} \frac{\partial \mathbf{z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = \sum_{i'} \mathbf{K}_{k,j} \sum_{\ell'=0}^{d-1} \mathbf{W}_{\ell',j}^2 \mathbb{E} [\mathbf{u}_{i',\ell'}^2] \mathbb{E} [\mathbf{u}_{i'-k,j}^2].$$

Since we know that the summation piece over ℓ' is always non-zero due to the even exponents on the \mathbf{u} terms and at least one of $\mathbf{W}_{\ell',j} \neq 0$. Therefore, if we set

$$\mathbb{E} \left[\frac{\partial \overline{\mathcal{L}}}{\partial \mathbf{K}_{k,j}} \right] = 0,$$

we get for all k , $\mathbf{K}_{k,j} = 0$. In other words,

$$\mathbf{K} = \mathbf{0}^{N \times d}.$$

Finally, via [Lemma D.51](#) we have for all ℓ :

$$\begin{aligned} \mathbb{E} \left[\frac{\partial \overline{\mathcal{L}}}{\partial \mathbf{W}_{\ell,j}} \right] &= \sum_{i',j'} \mathbb{E} \left[\frac{\partial \overline{\mathcal{L}}_{i',j'}}{\partial \mathbf{W}_{\ell,j}} \right] \\ &= \sum_{i'} \mathbb{E} \left[\frac{\partial \overline{\mathcal{L}}_{i',j}}{\partial \mathbf{W}_{\ell,j}} \right] \end{aligned}$$

From [Lemma D.57](#) and [Lemma D.60](#) we get

$$\begin{aligned} \sum_{i'} \mathbb{E} \left[\frac{\partial \overline{\mathcal{L}}_{i',j}}{\partial \mathbf{W}_{\ell,j}} \right] &= \sum_{i'} \mathbf{w}_{\ell,j} \sum_{k'=0}^{i'} \mathbf{K}_{k',j}^2 \mathbb{E} [\mathbf{u}_{i'-k',j}^2 \mathbf{u}_{i',\ell}^2] + \left(\mathbf{B}_{i',j}^{(2)} \right)^2 \mathbf{w}_{\ell,j} \mathbb{E} [\mathbf{u}_{i',\ell}^2] - \mathbf{B}_{i',j}^2 \overline{\mathbf{W}}_{\ell,j} \mathbb{E} [\mathbf{u}_{i',\ell}^2] \\ &= \mathbf{w}_{\ell,j} \sum_{i'} \left(\mathbf{B}_{i',j}^{(2)} \right)^2 \mathbb{E} [\mathbf{u}_{i',\ell}^2] - \overline{\mathbf{W}}_{\ell,j} \sum_{i'} \mathbf{B}_{i',j}^{(2)} \mathbb{E} [\mathbf{u}_{i',\ell}^2]. \end{aligned}$$

We can drop the summation with \mathbf{K} in it as we know $\mathbf{K} = \mathbf{0}$ Recall that from (57), $\mathbf{B}_{i,j}^{(2)} = b_j$. Then we can rewrite the above as

$$\mathbb{E} \left[\frac{\partial \overline{\mathcal{L}}}{\partial \mathbf{W}_{\ell,j}} \right] = b_j \left(\sum_{i'} \mathbb{E} [\mathbf{u}_{i',\ell}^2] \right) (\mathbf{w}_{\ell,j} b_j - \overline{\mathbf{W}}_{\ell,j})$$

we know from [Assumption D.45](#) that the first summation will always be non-zero since there's an even exponent on the \mathbf{u} variable and we know that b_j is non-zero. Therefore setting

$$\mathbb{E} \left[\frac{\partial \overline{\mathcal{L}}}{\partial \mathbf{W}_{\ell,j}} \right] = 0$$

tells us that

$$\mathbf{w}_{\ell,j} = \frac{\overline{\mathbf{W}}_{\ell,j}}{b_j}.$$

Given the above value for $\mathbf{W}_{\ell,j}$ and recall we have $\mathbf{K} = \mathbf{0}$ and $\mathbf{B}^{(2)} = (b_0 b_1 \dots b_{d-1})$ where each b_j is a column vector comprised of all b_j values. Therefore, when we take $\text{BASECONV}(\mathbf{u})$ we get

$$\begin{aligned} \text{BASECONV}(\mathbf{u}) &= (\mathbf{u}\mathbf{W}) \odot \left(\mathbf{0}^{N \times d} * \mathbf{u} + \mathbf{B}^{(2)} \right) \\ &= (\mathbf{u}\mathbf{W}) \odot \left(\mathbf{B}^{(2)} \right) \end{aligned}$$

We can rewrite $\mathbf{B}^{(2)}$ as

$$\mathbf{B}^{(2)} = (\mathbf{1}^{N \times d}) \begin{pmatrix} b_0 & 0 & \dots & 0 \\ 0 & b_1 & \dots & 0 \\ & & \ddots & \\ 0 & 0 & \dots & b_{d-1} \end{pmatrix}.$$

Lets call this diagonal matrix on the right, \mathbf{D} . Then note that

$$\mathbf{W} = \overline{\mathbf{W}}\mathbf{D}^{-1}.$$

Therefore, we have

$$\begin{aligned} \text{BASECONV}(\mathbf{u}) &= \mathbf{u}\overline{\mathbf{W}}\mathbf{D}^{-1} \odot \mathbf{1}^{N \times d}\mathbf{D} \\ &= \mathbf{u}\overline{\mathbf{W}} \odot \mathbf{1}^{N \times d} \\ &= \mathbf{u}\overline{\mathbf{W}}, \end{aligned}$$

as desired. In the above the second inequality follows since \mathbf{D} is a diagonal matrix. \square

We now revisit the importance of [Assumption D.55](#). Specifically, the following definition is a stronger version of the complement of [Assumption D.55](#). The following essentially states that there are many ways to get the expected gradients of the loss function to be 0, though this doesn't imply that we have learned the exact solution, as we recover in [theorem D.59](#).

Definition D.63. Define (assumption)^G to be

- $\mathbf{B}^{(2)} = \mathbf{0}$
- For all j , either
 - (i) $\mathbf{W}_{:,j} = \mathbf{0}$ and $\mathbf{K}_{0,j} = 0$
 - (ii) $\mathbf{K}_{:,j} = \mathbf{0}$ and $\mathbf{W}_{j,j} = 0$

The following theorem is to emphasize, there are many ways to get expected value of the gradients of the loss function to be 0.

Theorem D.64. Let θ^* satisfy [Definition D.63](#) Then $\mathbb{E}\nabla_{\theta}\overline{L}|_{\theta \leftarrow \theta^*} = 0$ when $f(\mathbf{u}) = \mathbf{u} \odot \mathbf{u}$

Proof. Next, via [Lemma D.57](#) and [Lemma D.58](#) when $k > 0$ we have

$$\mathbb{E} \left[\frac{\partial \overline{L}}{\partial \mathbf{K}_{k,j}} \right] = \mathbf{K}_{k,j} \sum_{i'} \sum_{\ell'=0}^{d-1} \mathbf{W}_{\ell',j}^2 \mathbb{E} [\mathbf{u}_{i',\ell'}^2] \mathbb{E} [\mathbf{u}_{i'-k,j}^2].$$

This expected value goes to zero since every column j either $\mathbf{K}_{:,j}$ or $\mathbf{W}_{:,j}$ is $\mathbf{0}$.

When $k = 0$ we have

$$\mathbb{E} \left[\frac{\partial \overline{L}}{\partial \mathbf{K}_{0,j}} \right] = \mathbf{W}_{j,j} (\mathbf{K}_{0,j} \mathbf{W}_{j,j} - 1) \left(\sum_{i'} \mathbb{E} [\mathbf{u}_{i',j}^4] \right).$$

This expected value goes to zero since in both [D.63](#) and [D.63](#) $\mathbf{K}_{0,j} = \mathbf{W}_{j,j} = 0$.

Next we have for all ℓ, j where $\ell \neq j$,

$$\mathbb{E} \left[\frac{\partial \overline{L}}{\partial \mathbf{W}_{\ell,j}} \right] = \mathbf{W}_{\ell,j} \sum_{i'} \mathbb{E} [\mathbf{u}_{i',\ell}^2] \left(\mathbf{K}_{0,j}^2 \mathbb{E} [\mathbf{u}_{i',j}^2] + \mathbf{B}_{i',j}^{(2)} \right)$$

3834 This expected value goes to zero since column j either $\mathbf{K}_{:,j}$ or $\mathbf{W}_{:,j}$ is $\mathbf{0}$.

3835 Then when $\ell = j$ we have

$$3837 \mathbb{E} \left[\frac{\partial \bar{L}}{\partial \mathbf{W}_{j,j}} \right] = \sum_{i'} \mathbf{W}_{j,j} \mathbf{K}_{0,j}^2 \mathbb{E} [\mathbf{u}_{i',j}^4] - \mathbf{K}_{0,j} \mathbb{E} [\mathbf{u}_{\ell',j}^4]$$

3840 This expected value goes to zero since in both (i) and (ii), $\mathbf{K}_{0,j} = \mathbf{W}_{j,j} = \mathbf{0}$

3841 And we know that since $\mathbf{B}^{(2)}$ is all zeros, we don't need to consider the gradient of the loss function
3842 to it. \square

3843 What the above proves is that there are infinite instantiations of parameters such that the expected
3844 gradient loss is 0. However not that in [Definition D.63](#), for all j either $\mathbf{K}_{k,j}$ for $k \neq 0$ or $\mathbf{W}_{\ell,j}$
3845 for $\ell \neq j$ are unconstrained. In other words, we can set these values arbitrarily, which means we
3846 get $\bar{L} \rightarrow \inf$ but we still have the expected gradient loss to be $\mathbf{0}$. This shows that some form of
3847 [Assumption D.55](#) is necessary to prove [Theorem D.59](#).

3850 D.8 LEARNING ANY POSSIBLE FUNCTION

3851 We hope to be able to show that one layer of BASECONV will be able to exactly recover any one
3852 layer BASECONV, given some assumptions on the parameters. We will later explain why these
3853 assumptions are necessary for a meaningful result.

3854 **Lemma D.65.** *Fix i, j . Then we have*

$$3857 \mathbb{E} \left[\left((\mathbf{u}\bar{\mathbf{W}}) \odot (\bar{\mathbf{K}} * \mathbf{u} + \bar{\mathbf{B}}^{(2)}) \right)_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] = \bar{\mathbf{B}}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'}^2] \bar{\mathbf{W}}_{\ell',j} \mathbf{W}_{\ell',j}.$$

3859 For $0 \leq k \leq i$, we have

$$3862 \mathbb{E} \left[\left((\mathbf{u}\bar{\mathbf{W}}) \odot (\bar{\mathbf{K}} * \mathbf{u} + \bar{\mathbf{B}}^{(2)}) \right)_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = \bar{\mathbf{K}}_{k',j} \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'}^2] [\mathbf{u}_{i-k,j}^2] \bar{\mathbf{W}}_{\ell',j} \mathbf{W}_{\ell',j}$$

3864 Next for all $0 \leq \ell < d$,

$$3867 \mathbb{E} \left[\left((\mathbf{u}\bar{\mathbf{W}}) \odot (\bar{\mathbf{K}} * \mathbf{u} + \bar{\mathbf{B}}^{(2)}) \right)_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] = \bar{\mathbf{W}}_{\ell,j} \sum_{k'=0}^i \mathbb{E} [\mathbf{u}_{i,\ell}^2] \mathbb{E} [\mathbf{u}_{i-k',j}^2] \bar{\mathbf{K}}_{k',j} \mathbf{K}_{k',j} + \bar{\mathbf{B}}_{i,j}^{(2)} \bar{\mathbf{W}}_{\ell,j} \mathbb{E} [\mathbf{u}_{i,\ell}^2]$$

3869 *Proof.* Lets begin by looking at

$$3872 \mathbb{E} \left[\left((\mathbf{u}\bar{\mathbf{W}}) \odot (\bar{\mathbf{K}} * \mathbf{u} + \bar{\mathbf{B}}^{(2)}) \right)_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right].$$

3874 From [lemma D.52](#) and [lemma D.54](#) we have

$$3875 \mathbb{E} \left[\left((\mathbf{u}\bar{\mathbf{W}}) \odot (\bar{\mathbf{K}} * \mathbf{u} + \bar{\mathbf{B}}^{(2)}) \right)_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] = \mathbb{E} \left[\left((\mathbf{u}\bar{\mathbf{W}}) \odot (\bar{\mathbf{K}} * \mathbf{u} + \bar{\mathbf{B}}^{(2)}) \right)_{i,j} \cdot \langle (\mathbf{u}_{i,:}, \mathbf{W}_{:,j}) \rangle \right]$$

3877 This expands to

$$3880 \sum_{\ell'=0}^{d-1} \sum_{\ell''=0}^{d-1} \sum_{k=0}^i \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell''} \mathbf{u}_{i-k,j}] \bar{\mathbf{W}}_{\ell',j} \mathbf{W}_{\ell'',j} \mathbf{K}_{k,j} + \bar{\mathbf{B}}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \sum_{\ell''=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell''}] \bar{\mathbf{W}}_{\ell',j} \mathbf{W}_{\ell'',j}.$$

3882 We simplify the above using [Assumption D.45](#). The first summation goes to 0 as we're taking the
3883 expectation of an odd number of \mathbf{u} 's. The second summation will be non-zero if and only if $\ell' = \ell''$,
3884 as that gives us a squared value on the \mathbf{u} variable. Therefore we get

$$3886 \mathbb{E} \left[\left((\mathbf{u}\bar{\mathbf{W}}) \odot (\bar{\mathbf{K}} * \mathbf{u} + \bar{\mathbf{B}}^{(2)}) \right)_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{B}^{(2)}_{i,j}} \right] = \bar{\mathbf{B}}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'}^2] \bar{\mathbf{W}}_{\ell',j} \mathbf{W}_{\ell',j}$$

3888 as desired.

3889 Next lets consider, for all $0 \leq k \leq i$,

$$3891 \mathbb{E} \left[\left((\mathbf{u}\overline{\mathbf{W}}) \odot (\overline{\mathbf{K}} * \mathbf{u} + \overline{\mathbf{B}^{(2)}}) \right)_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right].$$

3894 Recall we are only considering these values of k as we know from [Lemma D.54](#) that for $k > i$,
3895 the partial derivative piece is always 0. We can simplify this to the following via [lemma D.52](#) and
3896 [lemma D.54](#),

$$3897 \mathbb{E} \left[\left((\mathbf{u}\overline{\mathbf{W}}) \odot (\overline{\mathbf{K}} * \mathbf{u} + \overline{\mathbf{B}^{(2)}}) \right)_{i,j} \langle (\mathbf{u}_{i,:}, \mathbf{W}_{:, \ell}) \mathbf{u}_{i-k,j} \rangle \right]$$

3899 This expands to

$$3901 \sum_{\ell'=0}^{d-1} \sum_{\ell''=0}^{d-1} \sum_{k'=0}^i \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell''} \mathbf{u}_{i-k',j} \mathbf{u}_{i-k,j}] \overline{\mathbf{W}}_{\ell',j} \overline{\mathbf{W}}_{\ell'',j} \overline{\mathbf{K}}_{k',j} + \overline{\mathbf{B}}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \sum_{\ell''=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell''} \mathbf{u}_{i-k,j}] \overline{\mathbf{W}}_{\ell',j} \overline{\mathbf{W}}_{\ell'',j}$$

3904 We simplify the above using [Assumption D.45](#). The first summation will be non-zero if and only if
3905 the exponents on the \mathbf{u} 's is even. To get this we require $\ell' = \ell''$ and $k' = k$. The second summation
3906 will always be 0 as we have an odd exponent on the \mathbf{u} 's. Therefore we get

$$3908 \mathbb{E} \left[\left((\mathbf{u}\overline{\mathbf{W}}) \odot (\overline{\mathbf{K}} * \mathbf{u} + \overline{\mathbf{B}^{(2)}}) \right)_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{K}_{k,j}} \right] = \overline{\mathbf{K}}_{k,j} \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'}^2] [\mathbf{u}_{i-k,j}^2] \overline{\mathbf{W}}_{\ell',j} \overline{\mathbf{W}}_{\ell',j}$$

3911 as desired.

3912 Finally, lets consider

$$3914 \mathbb{E} \left[\left((\mathbf{u}\overline{\mathbf{W}}) \odot (\overline{\mathbf{K}} * \mathbf{u} + \overline{\mathbf{B}^{(2)}}) \right)_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right]$$

3916 From [lemma D.52](#) and [lemma D.54](#) we get

$$3918 \mathbb{E} \left[\left((\mathbf{u}\overline{\mathbf{W}}) \odot (\overline{\mathbf{K}} * \mathbf{u} + \overline{\mathbf{B}^{(2)}}) \right)_{i,j} \left((\mathbf{K}_{:,j} * \mathbf{u}_{:,j}) [i] + \mathbf{B}_{i,j}^{(2)} \right) \mathbf{u}_{i,\ell} \right]$$

3920 This expands to

$$3922 \sum_{\ell'=0}^{d-1} \sum_{k'=0}^i \sum_{k''=0}^i \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell} \mathbf{u}_{i-k',j} \mathbf{u}_{i-k'',j}] \overline{\mathbf{W}}_{\ell',j} \overline{\mathbf{K}}_{k',j} \overline{\mathbf{K}}_{k'',j} + \overline{\mathbf{B}}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \sum_{k''=0}^i \mathbb{E} [\mathbf{u}_{i-k'',j} \mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell'}] \overline{\mathbf{W}}_{\ell',j} \overline{\mathbf{K}}_{k'',j}$$

$$3924 + \overline{\mathbf{B}}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \sum_{k'=0}^i \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i-k',j} \mathbf{u}_{i,\ell}] \overline{\mathbf{W}}_{\ell',j} \overline{\mathbf{K}}_{k',j} + \overline{\mathbf{B}}_{i,j}^{(2)} \overline{\mathbf{B}}_{i,j}^{(2)} \sum_{\ell'=0}^{d-1} \mathbb{E} [\mathbf{u}_{i,\ell'} \mathbf{u}_{i,\ell}] \overline{\mathbf{W}}_{\ell',j}.$$

3928 We simplify the above using [Assumption D.45](#). The first summation will be non-zero as we get
3929 squared exponents on the \mathbf{u} 's if and only if $\ell' = \ell$ and $k' = k''$. The second and third summation will
3930 always be 0 since we're taking the expectation of an odd number of \mathbf{u} 's. The fourth summation will
3931 be non-zero if and only if $\ell' = \ell$. Therefore we get,

$$3933 \mathbb{E} \left[\left((\mathbf{u}\overline{\mathbf{W}}) \odot (\overline{\mathbf{K}} * \mathbf{u} + \overline{\mathbf{B}^{(2)}}) \right)_{i,j} \frac{\partial \mathbf{Z}_{i,j}}{\partial \mathbf{W}_{\ell,j}} \right] = \overline{\mathbf{W}}_{\ell,j} \sum_{k'=0}^i \mathbb{E} [\mathbf{u}_{i,\ell}^2] \mathbb{E} [\mathbf{u}_{i-k',j}^2] \overline{\mathbf{K}}_{k',j} \overline{\mathbf{K}}_{k',j} + \overline{\mathbf{B}}_{i,j}^{(2)} \overline{\mathbf{B}}_{i,j}^{(2)} \overline{\mathbf{W}}_{\ell,j} \mathbb{E} [\mathbf{u}_{i,\ell}^2]$$

3936 as desired. \square