

---

# MONOIDREDUCE: AN ALGEBRAIC FRAMEWORK FOR MEMORY-EFFICIENT NEURAL NETWORK LAYERS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Recent advances in memory-efficient neural network layers, such as FlashAttention, often appear as specialized engineering solutions but share a common mathematical structure. We show that many of these kernels can be understood as folds over commutative monoids, a perspective that unifies MapReduce-style computation with modern deep learning optimizations. Building on this, we introduce the Local Gradient Theorem, which provides a sufficient condition under which gradients of monoidal folds can be computed locally from the final output and individual inputs, enabling efficient backward passes. We demonstrate that attention, cross entropy, and two-layer MLPs all admit such monoid structures, recovering known memory-efficient kernels and extending the framework to new settings. This algebraic perspective offers a principled foundation for systematically designing memory- and cache-efficient layers, rather than discovering them in an ad-hoc manner.

## 1 INTRODUCTION

Training large neural networks is increasingly constrained by memory bandwidth and capacity rather than arithmetic throughput. In particular, attention layers and large-vocabulary classification heads require storing intermediate activations whose size grows quadratically or linearly with sequence length or vocabulary size. Recent work has shown that carefully reordering computation can eliminate the need to materialize these large intermediates. FlashAttention (Dao et al., 2022) introduced an IO-aware tiling strategy that computes exact attention in linear memory, later extended in FlashAttention-2 (Dao, 2023) for greater parallelism. In classification, Cut Cross-Entropy (CCE) (Wijmans et al., 2024) reorganizes the loss computation to avoid constructing full logit matrices. Liger kernels (Hsu et al., 2024) take a similar philosophy to operator fusion, providing efficient Triton implementations of common primitives. These works illustrate a pattern: memory efficiency can be achieved when a kernel can be expressed as a reduction.

A complementary line of work reduces memory through activation recomputation. General strategies such as gradient checkpointing (Chen et al., 2016; Griewank and Walther, 2000) selectively discard activations and recompute them during the backward pass, while reversible networks (Gomez et al., 2017) reconstruct hidden states exactly from later layers. IO-aware kernels such as FlashAttention (Dao et al., 2022; Dao, 2023) and Cut Cross-Entropy (Wijmans et al., 2024) can be viewed as specialized instances of the former: they deliberately recompute partial results in an online fashion to avoid materializing the full intermediate tensor in memory.

Our work generalizes this observation: we identify a broad class of *monoidal reductions* where such structured recomputation is always possible under a simple local gradient condition, and propose **MonoidReduce**, an abstraction that unifies FlashAttention, Cut Cross-Entropy, and related fused kernels, by observing that many of these algorithms can be written as *folds* over commutative monoids, enabling streaming and tiling. We prove a *Local Gradient Theorem* that characterizes when such folds admit recompute-friendly backwards passes. This provides a simple algebraic condition under which kernels can be made memory-efficient. A Proof of Concept implementation of this abstraction can be found in Appendix B.

Finally, our framework connects modern ML kernel design to classical systems and parallel algorithms. The fold view is closely related to MapReduce (Dean and Ghemawat, 2008) and prefix-sum

based reductions (Blelloch, 1990), providing a bridge between deep learning practice and foundational results in parallel computation. This is not novel, with recent work bridging this gap (Rush et al., 2024b;a). However, these works do not address the memory-efficient backwards passes that follow from the *Local Gradient Theorem*.

### 1.1 WHAT ARE COMMUTATIVE MONOIDS?

A monoid is a tuple  $(T, \odot, \text{id})$ , where  $T$  is a type,  $\odot : T \rightarrow T \rightarrow T$  a binary operation on  $T$ , and  $\text{id}$  a term of type  $T$ . The binary operation must be associative, i.e.  $a \odot (b \odot c) = (a \odot b) \odot c$ , and  $\text{id}$  must be an identity element, i.e.  $\text{id} \odot a = a = a \odot \text{id}$ . Commutative monoids are monoids for which the binary operation is also commutative, that is,  $a \odot b = b \odot a$ .

There are many instances of commutative monoids in modern deep learning:

**Sum** Summation over real numbers:

$$\begin{aligned} T &= \mathbb{R} \\ a \odot b &= a + b \\ \text{id} &= 0 \end{aligned}$$

**WSum** Weighted sums of vectors.

$$\begin{aligned} T &= \{v : \mathbb{R}^D, w : \mathbb{R}_{\geq 0}\} \\ a \odot b &= c \\ \text{where } c_w &= a_w + b_w \\ c_v &= \begin{cases} a_v \frac{a_w}{c_w} + b_v \frac{b_w}{c_w}, & c_w > 0 \\ 0, & \text{otherwise} \end{cases} \\ \text{id} &= \{v : 0, w : 0\} \end{aligned}$$

**LogWSum** Weighted sums of vectors with weights in logspace<sup>1</sup>:

$$\begin{aligned} T &= \{v : \mathbb{R}^D, z : \mathbb{R}\} \\ a \odot b &= c \\ \text{where } c_z &= \ln(\exp(a_z) + \exp(b_z)) \\ c_v &= \begin{cases} a_v \exp(a_z - c_z) + b_v \exp(b_z - c_z), & c_z > -\infty \\ 0, & \text{otherwise} \end{cases} \\ \text{id} &= \{v : 0, z : -\infty\} \end{aligned}$$

Given a commutative monoid  $T$ , we can derive a function *fold* that takes as input a sequence of  $T$ s and returns the product over all elements in the sequence:  $\text{fold} : [T] \rightarrow T = H \mapsto \bigodot_i H_i$ .

## 2 MONOIDREDUCE

Commutative monoidal folds bring several advantages to large-scale machine learning. First, because the fold operation can be applied incrementally, there is no need to materialize all intermediate values in memory during the forward pass. This makes it possible to stream or distribute data through the computation, using GPU memory more efficiently, much like how MapReduce (Dean and Ghemawat, 2008) avoids keeping the entire dataset in memory by processing records in batches. Second, since partial results of a fold can be combined in any order, they can be computed independently across different devices (in the distributed setting) or thread blocks (in the kernel setting) and later merged, analogous to the “reduce” stage in MapReduce. Finally, if the derivative of the monoidal product ( $\frac{d}{dx} x \odot y$ ) can be expressed in terms of  $x \odot y$  and  $x$ , then the gradient of the entire fold inherits this property. This means gradients can be computed from the final aggregated value and *local* recomputation, enabling parallelization of the backward pass as well. This last observation is the main novelty presented in this paper, which we formalize as the *Local Gradient Theorem*,

<sup>1</sup>This monoid applies to attention:  $\text{softmax}(A)V = x_0 \odot x_1 \odot \dots$ , where  $x_i = \{v = V_i, z = A_i\}$

**Theorem 2.1** (Local Gradient Theorem). *Let  $T$  be a finite-dimensional vector space and  $(T, \odot, \mathbf{id})$  a commutative monoid with a derivative<sup>2</sup>  $\frac{d x \odot y}{d x}$  that can be expressed as a function  $(D)$  of  $x \odot y$  and  $x$ :*

$$\frac{d x \odot y}{d x} = D(x \odot y, x) \quad (1)$$

We have that for any sequence  $H : [T]$ , the following holds:

$$\frac{d \text{fold}(H)}{d H_i} = D(\text{fold}(H), H_i)$$

*Proof.* Let  $P = \text{fold}(H) = \odot_i H_i$ , for any partition of  $H$  into  $H^1, H^2$ , with corresponding partial products  $P_1 = \text{fold}(H^1) = \odot_i H_i^1$  and  $P_2 = \text{fold}(H^2) = \odot_i H_i^2$ .

$$\begin{aligned} P &= P_1 \odot P_2 = P_2 \odot P_1 && \text{(By commutativity and associativity)} \\ \frac{d P}{d P_1} &= \frac{d P_1 \odot P_2}{d P_1} \\ &= D(P_1 \odot P_2, P_1) && \text{(By the condition from eq. 1)} \\ &= D(P, P_1) \end{aligned}$$

This extends trivially to  $P_1 = H_i$ , i.e.,  $H^1 = [H_i]$ :

$$\begin{aligned} \frac{d P}{d H_i} &= \frac{d P_1}{d H_i} \cdot D(P, P_1) \\ &= D(P, P_1) && \left( \frac{d P_1}{d H_i} = g \mapsto g \right) \\ &= D(P, H_i) \end{aligned}$$

□

**Corollary 2.1.1.** *Consider the scenario where  $H_{ij} = f(A_i, B_j)$  and let  $P$  be the vector of final products ( $P_i = \text{fold}(H_i) = \odot_j H_{ij}$ ). If we satisfy eq. 1, then we have that:*

$$\begin{aligned} \frac{d P}{d A_i} &= \sum_j \frac{d H_{ij}}{d A_i} \cdot D(P_i, H_{ij}) \cdot \frac{d P}{d P_i} \\ \frac{d P}{d B_j} &= \sum_i \frac{d H_{ij}}{d B_j} \cdot D(P_i, H_{ij}) \cdot \frac{d P}{d P_i} \end{aligned}$$

And more generally that for any partition of  $A$  into  $A^1, A^2, \dots, A^m$  and  $B$  into  $B^1, B^2, \dots, B^n$ , with corresponding partitions of  $H$  into  $H^{11}, H^{12}, \dots, H^{mn}$  and partial product vectors  $P^{ij}$ , where  $P_k^{ij} = \odot_l H_{kl}^{ij}$ , and let  $P^i$  denote the vector of products  $P_k^i = \odot_j P_k^{ij}$ :

$$\begin{aligned} \frac{d P}{d A^i} &= \sum_j \frac{d P^{ij}}{d A^i} \cdot D(P^i, P^{ij}) \cdot \frac{d P}{d P^i} \\ \frac{d P}{d B^j} &= \sum_i \frac{d P^{ij}}{d B^j} \cdot D(P^i, P^{ij}) \cdot \frac{d P}{d P^i} \end{aligned}$$

<sup>2</sup>The mathematical convention used here is that  $\frac{d A}{d B}$  denotes the linear map that transforms gradients of  $A$  into gradients of  $B$  (backward / pullback Jacobian), inspired by Conal Elliot's work on Automatic Differentiation (Elliott, 2009). Function composition  $\cdot$  denotes sequential application of these linear maps, so the chain rule reads  $\frac{d c}{d a} = \frac{d b}{d a} \cdot \frac{d c}{d b}$ , corresponding to  $(f \cdot g)(x) = f(g(x))$ .  $\langle a, b \rangle$  denotes dot product of vectors.

Remark:  $\frac{dP}{dP_i}$  and  $\frac{dP}{dP^i}$  are simple diagonal projections.

The implication of Corollary 2.1.1 is that gradients can be computed based only on local recomputation and final products, *not* the full  $H$ -matrix, which in turn facilitates highly parallel and cache-efficient kernels for backward passes.

In Appendix B we show a general proof of concept implementation of the MonoidReduce abstraction (listing 2), and an implementation of attention that uses this abstraction (listing 3), where forward and backward passes rely on user-supplied functions that map naturally to corollary 2.1.1. Table 1 summarizes these functions and their theoretical correspondence.

User Function	Pass	Description
<code>proj_fold(A, B)</code>	Forward	Computes the partial fold vector $P^{ij}$ for slices $A^i, B^j$ of $A$ and $B$ .
<code>proj_fold_bwd(A, B, P, gP)</code>	Backward	Computes local gradients $\frac{dP^{ij}}{dA^i} \cdot D(P^i, P^{ij})$ and $\frac{dP^{ij}}{dB^j} \cdot D(P^i, P^{ij})$ .
<code>binary_reduce(P1, P2)</code>	Forward	The (vectorized) monoid operation $\odot$ .
<code>init(A, B)</code>	Forward	Produces an appropriately shaped (monoidal) identity-valued tensor.
<code>islice(A), jslice(B)</code>	Both	Produces partitions of the input, output, and gradients corresponding to $A^i, P^i, B^j$ , etc.

Table 1: User-supplied functions for batched MonoidReduce and their correspondence to corollary 2.1.1

**Forward pass.** Each slice pair  $(A^i, B^j)$  is processed via `proj_fold` to produce partial results  $P^{ij}$ , which are then combined with (vectorized) `binary_reduce` to obtain the slice-level outputs  $P^i$ . This is embarrassingly parallel over  $i$ , but parallelizing over  $j$  as well requires parallel reduction of the monoidal product  $\odot$ . A visualization of this can be seen in Figure 1.

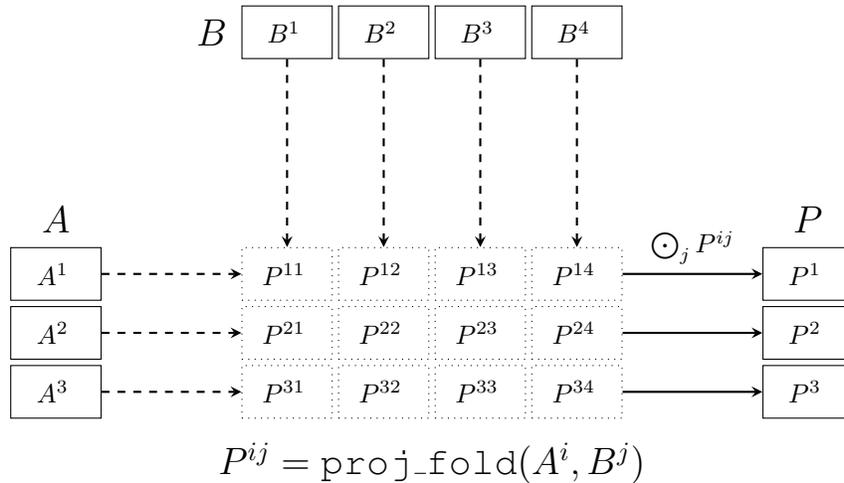


Figure 1: A visualization of the MonoidReduce forward pass.

**Backward pass.** Gradients for each slice are computed using `proj_fold_bwd` for each pair of slices, and the results are added together across slices to form the full gradients for  $A$  and  $B$ . The local gradient theorem (2.1) ensures that this is sound. Note that the summation of partial gradients poses a potential communication overhead. A visualization of this can be seen in Figure 2.

Corollary 2.1.1 also provides a general method for computing the Jacobian for  $A$  during a fold over  $B$ , limited by how efficiently the Jacobian sum can be represented. In the case where  $T$  is essentially a scalar (e.g. where  $a \odot b = \log\text{addexp}(a, b)$  i.e.  $\text{fold}(H) = \log\text{sumexp}(H)$ ) this is manageable as the resulting Jacobian has the same shape as  $A$ , but for more complex monoids the

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

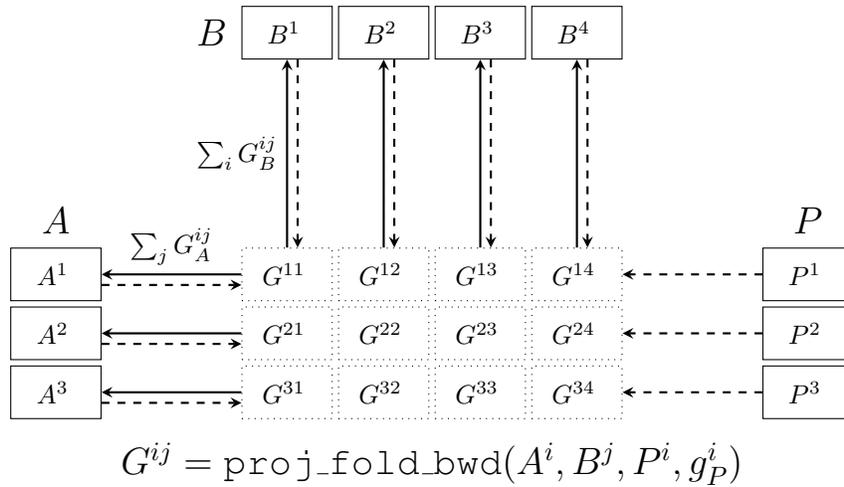


Figure 2: A visualization of the MonoidReduce backward pass.

cost of keeping the Jacobian in memory might be prohibitively expensive. When possible, such an approach simplifies the backward pass further, enabling us to apply the precomputed Jacobian to  $\frac{d_j \mathcal{L}}{d P}$  to get the gradients for  $A$ , while separately computing the  $B$ -gradient by parallelizing over slices of  $B$ . As done in Hsu et al. (2024); Wijmans et al. (2024).

## 2.1 EXAMPLE: ATTENTION

The benefit of not materializing all partial values is most apparent in cases where the fold is applied to the results of matrix multiplications; one such instance is attention:

$$\begin{aligned}
 Q &: \mathbb{R}^{M \times F} \\
 K &: \mathbb{R}^{N \times F} \\
 V &: \mathbb{R}^{N \times D} \\
 \text{attention}(Q, K, V) &= \text{softmax}(QK^\top)V \\
 \text{attention}(Q, K, V)_i &= \left( \bigodot_j f(A_i, B_j) \right)_v
 \end{aligned}$$

where

$$\begin{aligned}
 A_i &= Q_i \\
 B_j &= \{k : K_j, v : V_j\} \\
 f(a, b) &= \{z : \langle a, b_k \rangle, v : b_v\} \\
 a \odot b &= \left\{ \begin{array}{l} z : \ln(e^{a_z} + e^{b_z}) \\ v : a_v e^{a_z - z} + b_v e^{b_z - z} \end{array} \right\}
 \end{aligned}$$

Here, we can think of  $f$  as realizing an  $M$  by  $N$  matrix of  $T$ -values. Each  $T$ -value consists of a log-space weight  $z$ , and a value-vector  $v$ . The monoidal product  $a \odot b$  results in a new  $T$ , with a weight corresponding to the log-space sum of  $a$ 's and  $b$ 's weight, and a value-vector corresponding to the weighted average of  $a$ 's and  $b$ 's value-vectors. Since the result is a fold over the  $N$ -dimension, the full matrix does not have to be materialized. We can compute appropriate chunks of the matrix, fold the chunk, and add to a running total, only realizing the  $M$  normalizing factors, and the  $M \times D$ -sized weighted sum.

We also have that condition 1 holds:

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

$$\begin{aligned}
c &= a \odot b \\
\frac{dc}{da} &= g \mapsto \begin{pmatrix} z: & (g_z + \langle g_v, a_v - c_v \rangle) \exp(a_z - c_z) \\ v: & g_v \exp(a_z - c_z) \end{pmatrix} \\
&= D(c, a)
\end{aligned}$$

Which due to Corollary 2.1.1 implies that the gradient w.r.t.  $Q$ ,  $K$  and  $V$  can be calculated based on the final product  $P$  and local (recomputed)  $T$ -values. Implementations of the necessary functions from Table 1 can be seen in Listing 1. Flash attention (Dao et al., 2022) can be seen as a highly optimized instance of this general approach.

Listing 1: Instances of functions from Table 1 for attention.

```

def init(a, b):
    aq, = a
    _, bv = b
    L, _ = aq.shape
    _, D = bv.shape
    z = aq.new_full((L, ), float("-inf"))
    v = bv.new_zeros((L, D))

def proj_fold(a, b):
    aq, = a
    bk, bv = b
    logits = aq @ bk.t()
    z = torch.logsumexp(logits, dim=1)
    v = (logits - z[:, None]).exp() @ bv
    return z, v

def binary_reduce(x, y):
    xz, xv = x
    yz, yv = y
    z = torch.logaddexp(xz, yz)
    v = (xz - z).exp()[:, None] * xv + (yz - z).exp()[:, None] * yv
    return z, v

def D(global_p, local_p, gp):
    pz, pv = global_p
    lpz, lpv = local_p
    gz, gv = gp
    lgz = (gz + (gv * (lpv - pv)).sum(1)) * (lpz - pz).exp()
    lgv = gv * (lpz - pz).exp()[:, None]
    return lgz, lgv

def proj_fold_bwd(a, b, p, gp):
    # Fusing this would be more efficient.
    # The general form is used for presentation purposes.
    local_p, local_vjp = torch.func.vjp(proj_fold, a, b)
    return local_vjp(D(p, local_p, gp))

```

In Appendix C, Table 2, we give further examples of this approach applied to cross entropy against class indices<sup>3</sup>, cross entropy between two distributions, and two-layer MLPs<sup>4</sup>, all of which satisfy the condition for the local gradient theorem (Table 3).

### 3 COMPLEXITY ANALYSIS

The time and space complexity of MonoidReduce compared to a naive implementation depends on the time and space complexity of `proj_fold`, `proj_fold_bwd`, and `binary_reduce`, and the slicing of  $A$  and  $B$ .

<sup>3</sup>Similar to both Liger Kernels and Cut Your Losses (Hsu et al., 2024; Wijmans et al., 2024)

<sup>4</sup>Similar to Tensor Parallelism in the multi-device setting (Shoeybi et al., 2019)

324 Consider inputs  $A : \mathcal{A}^L$  and  $B : \mathcal{B}^R$ . We partition  $A$  into  $l$  equal-sized slices  $A^1, \dots, A^l$  and  $B$   
 325 into  $r$  equal-sized slices  $B^1, \dots, B^r$ , let  $|\mathcal{A}|$ , let  $|\mathcal{B}|$ , and let  $|\mathcal{P}|$  be the size of individual values of  
 326 type  $\mathcal{A}$ ,  $\mathcal{B}$ , and  $\mathcal{P}$ , respectively. Furthermore assume time and space complexities are (bi)linear with  
 327 regards to  $m$  and  $n$ .<sup>5</sup>  
 328  
 329

	time	space	assumption
330 $\text{proj\_fold} : \mathcal{A}^m \times \mathcal{B}^n \rightarrow \mathcal{P}^m$	$T_f(m, n)$	$S_f(m, n)$	bilinear
331 $\text{proj\_fold\_bwd} : \mathcal{A}^m \times \mathcal{B}^n \times \mathcal{P}^m \times \mathcal{P}^m \rightarrow \mathcal{A}^m \times \mathcal{B}^n$	$T_b(m, n)$	$S_b(m, n)$	bilinear
332 $\text{binary\_reduce} : \mathcal{P}^m \times \mathcal{P}^m \rightarrow \mathcal{P}^m$	$T_p(m)$	$S_p(m)$	linear

### 3.1 FORWARD PASS

333 **Time Complexity.** The MonoidReduce forward pass uses  $l \times r$  calls to `proj_fold`, each taking  
 334 time  $T_f(\frac{L}{l}, \frac{R}{r})$ , with  $l \times (r - 1)$  reductions, each taking time  $T_p(\frac{L}{l})$ .  
 335

336 The total time complexity is therefore  $l \times r \times T_f(\frac{L}{l}, \frac{R}{r}) + l \times (r - 1) \times T_p(\frac{L}{l})$ , which by (bi)linearity  
 337 simplifies to:

$$338 T_f(L, R) + (r - 1) \times T_p(L) \quad (2)$$

339 MonoidReduce thus incurs an extra cost of  $(r - 1) \times T_p(L)$  compared to the full `proj_fold`.

340 The critical path (with parallel reduction) takes time  $\frac{T_f(L, R)}{lr} + \frac{T_p(L) \times \log r}{l}$ , using Brent's Theorem  
 341 (Brent, 1974) this translates to an  $N$ -worker parallel time complexity of:

$$342 \frac{T_f(L, R) + (r - 1) \times T_p(L)}{N} + \frac{T_f(L, R)}{lr} + \frac{T_p(L) \times \log r}{l} \quad (3)$$

343 **Space Complexity.** For any particular worker, we only need to realize  $P^{ij}$ , and aggregate it to the  
 344 running total. Computing  $P^{ij}$  uses  $S_f(\frac{L}{l}, \frac{R}{r})$  space, and aggregation uses  $S_p(\frac{L}{l})$  space. Assuming  
 345 bilinearity, the per worker space complexity is therefore:

$$346 \frac{S_f(L, R)}{lr} + \frac{S_p(L)}{l} \quad (4)$$

347 A significant improvement compared to the full `proj_fold`.

348 **Communication Complexity.** The computation of  $P^{ij}$  requires reading  $A^i$  and  $B^j$ , incurring a  
 349 cost of  $\frac{L}{l} \times |\mathcal{A}| + \frac{R}{r} \times |\mathcal{B}|$ . The parallel sum  $\odot_j P^{ij}$  incurs a cost of  $(r - 1) \times \frac{L}{l} \times |\mathcal{P}|$ . Since we  
 350 need to compute  $l \times r$   $P^{ij}$ -values, and  $l$  parallel sums, this results in a total communication cost of:

$$351 r \times L \times |\mathcal{A}| + l \times R \times |\mathcal{B}| + (r - 1) \times L \times |\mathcal{P}| \quad (5)$$

352 MonoidReduce thus incurs an extra (total) communication cost of  $(r - 1) \times L \times (|\mathcal{A}| + |\mathcal{P}|) + (l - 1) \times$   
 353  $R \times |\mathcal{B}|$ . This can be mitigated grouping tasks along the fold-axis, letting one worker be responsible  
 354 for  $W$  slices  $B^j$ , resulting in a lower communication cost of  $\frac{r - W}{W} \times L \times (|\mathcal{A}| + |\mathcal{P}|) + (l - 1) \times R \times |\mathcal{B}|$   
 355 at the cost of worse idealized parallel time. This has been applied in, for example, FlashAttention  
 356 (Dao, 2023).  
 357

### 3.2 BACKWARD PASS

358 **Time Complexity.** The analysis is similar to that of the forward pass, but instead of monoidal fold  
 359 over  $P^{ij}$ , we aggregate gradients  $\sum_i G_A^{ij}$  and  $\sum_j G_B^{ij}$ .  
 360  
 361

$$362 T_b(L, R) + L \times |\mathcal{A}| + R \times |\mathcal{B}| \quad (6)$$

363 <sup>5</sup>This is generally the case when `proj_fold` is dominated by matrix multiplications of  $\mathbb{R}^{m \times k}$ ,  $\mathbb{R}^{k \times n}$ -  
 364 shaped matrices, and it is the case for all examples in Table 2.  
 365  
 366  
 367  
 368  
 369  
 370  
 371  
 372  
 373  
 374  
 375  
 376  
 377

i.e., identical to the full `proj_fold_bwd` followed by accumulating the gradients.

The critical path (with parallel reduction) takes time  $\frac{T_b(L,R)}{lr} + \frac{L}{l} \times |\mathcal{A}| \times \log r + \frac{R}{r} \times |\mathcal{B}| \times \log l$ , using Brent’s Theorem, this translates to an  $N$ -worker parallel time complexity of:

$$\frac{T_b(L, R) + L \times |\mathcal{A}| + R \times |\mathcal{B}|}{N} + \frac{T_b(L, R)}{lr} + \frac{L}{l} \times |\mathcal{A}| \times \log r + \frac{R}{r} \times |\mathcal{B}| \times \log l \quad (7)$$

**Space Complexity.** The analysis is similar to that of the forward pass. For any particular worker, we need only compute the local backward pass and add the resulting partial gradients to the global gradients:

$$\frac{S_b(L, R)}{lr} + \frac{L}{l}|\mathcal{A}| + \frac{R}{r}|\mathcal{B}| \quad (8)$$

A significant improvement compared to the full `proj_fold_bwd`.

**Communication Complexity.** The computation of  $G^{ij}$  requires reading  $A^i$ ,  $B^j$ ,  $P^i$ , and  $g_P^i$ . Incurring a communication cost of  $\frac{L}{l} \times |\mathcal{A}| + \frac{R}{r}|\mathcal{B}| + 2 \times \frac{L}{l}|\mathcal{P}|$ .

The parallel sums incur a cost of  $(r - 1) \times \frac{L}{l} \times |\mathcal{A}|$  and  $(l - 1) \times \frac{R}{r} \times |\mathcal{B}|$ .

This results in a total communication cost of:

$$2 \times r \times L \times (|\mathcal{A}| + |\mathcal{P}|) + 2 \times l \times R \times |\mathcal{B}| \quad (9)$$

The `MonoidReduce` backward pass thus incurs an extra communication cost of  $2 \times (r - 1) \times L \times (|\mathcal{A}| + |\mathcal{P}|) + 2 \times (l - 1) \times R \times |\mathcal{B}|$ . A way to mitigate this is to let one worker be responsible for several (input) slices  $A^i$  or  $B^j$ , accumulating gradients along those slices locally. If a worker is responsible for  $W$   $B^j$  (input) slices, thus accumulating  $W$   $A$ -gradients locally before adding to the global gradient, this reduces the communication overhead to  $2 \times \frac{r-W}{W} \times L \times (|\mathcal{A}| + |\mathcal{P}|) + 2 \times (l - 1) \times R \times |\mathcal{B}|$ ; again at the cost of worse idealized parallel time (critical path).

## 4 LIMITATIONS

While `MonoidReduce` provides an abstraction covering a wide range of existing memory-efficient implementations, and a *sufficient* condition for memory-efficient implementations in general, it does not apply universally. The central assumptions are that

- The reduction operator forms a commutative monoid.
- The derivative of the monoidal product can be expressed using only the combined value and one operand.
- The tradeoff between recomputation (or communication) cost and memory-efficiency is favorable.

This excludes many common operations, such as (sequences of) matrix multiplications, which can be understood as associative, *but not commutative*, folds. Instances where the tradeoff between recomputation (or communication) cost and memory-efficiency is not favorable are also common: Two-layer MLPs where the dimension being folded over is comparable in size to the output dimension will not see large memory-savings, and pay for the meager savings by incurring extra communication and computation.

From a systems perspective, recomputation introduces extra floating-point operations and potential hardware bottlenecks. In particular, as exemplified in Section 3, our generalized `MonoidReduce` example in listing 2 requires parallel fold of the monoidal product in the forward pass, and parallel sum of partial gradients in the backward pass; without careful tiling and shared-memory aggregation, atomic contention and communication costs can erode the benefits of theoretical memory savings. Furthermore, care must be taken with regard to numerical stability when implementing `proj_fold`, `proj_fold_bwd`, and `binary_reduce`.

---

432 Finally, our current work focuses on the algebraic formulation and theoretical guarantees rather  
433 than empirical benchmarking. Demonstrating practical wall-clock gains requires optimized imple-  
434 mentations and careful IO analysis, as done in FlashAttention (Dao et al., 2022; Dao, 2023). As  
435 such, this paper should be considered mainly theoretical, highlighting the conceptual foundation for  
436 memory-efficient layers.

## 437 438 5 DISCUSSION 439

440 As previously stated, this paper is mainly theoretical. The working implementation given in Ap-  
441 pendix B, while runnable and sometimes competitive against regular pytorch implementations on  
442 CPU, would require substantial work to go from Proof of Concept to fully fledged competitive  
443 framework. It is neither parallel nor very efficient. As hinted at in Section 3, slicing, grouping  
444 axis, and grouping size ( $W$ ), all have significant impact on efficiency. A future framework would  
445 require automatic tuning of these parameters, as they are crucial for finding a performant tradeoff  
446 between increased communication and decreased memory pressure. However, given such a frame-  
447 work, it would enable faster iteration and innovation in memory and cache-efficient layers, designed  
448 by design, rather than by accident.

449 Also note that MonoidReduce can produce values that are used for aggregation but not passed on-  
450 ward in the computational graph, such as the normalizing factor  $z$  in attention. This can enable more  
451 aggressive fusion of the `proj_fold_bwd` function and reduce communication cost. In a similar  
452 vein, if intermediate values used in the backward pass are shared within rows, they can be precom-  
453 puted, decreasing both compute and communication costs. The framework as presented does not  
454 adequately address such optimizations.

## 455 456 6 FUTURE WORK 457

458 The most straightforward avenue for future work would be a working and efficient implementation of  
459 the proposed framework, e.g., a GPU-aware implementation of listing 2 integrated into pytorch, with  
460 support for automatic tuning. This would be a significant undertaking. While probably not matching  
461 the peak performance of highly optimized, hand-tuned, single-purpose kernels, a well-engineered  
462 framework could offer competitive performance for many layers and significantly accelerate the  
463 development of new memory-efficient architectures.

464 One such avenue is MLPs: Let  $X, P, Q$  be matrices of shape  $B \times D$ ,  $K \times D$ , and  $K \times D$ , respec-  
465 tively, where  $B$  and  $K$  are very large, but  $D$  is small. Naive implementation would result in space  
466 complexity  $BD + 2KD + BK$ , with  $BK$  being the dominant factor. MonoidReduce can remove the  
467  $BK$ -term, at a cost of extra FLOPs during recomputation. The total FLOPs (forward and backward)  
468 for the naive implemntation is  $12BKD$ , whereas MonoidReduce increases this to  $14BKD$ . i.e.  
469  $\approx 17\%$  increase in FLOPs. If we let  $K = B = 16384$ , and  $D = 128$ , this would result in a memory  
470 requirement of  $\approx 2\%$  of the naive implementation, making MLPs for large batches ( $B$ ) and large  
471 hidden dimension ( $K$ ), but small input and output dimension ( $D$ ), much more feasible.

## 472 473 7 LLM USAGE 474

475 LLMs have been used in this work to flesh out drafts of the Introduction, Limitations, and Abstract,  
476 which were later edited. LLMs have also been used for general feedback and criticism.

## 477 478 REFERENCES 479

- 480 Blleloch, G. E. (1990). Prefix sums and their applications.  
481  
482 Brent, R. P. (1974). The parallel evaluation of general arithmetic expressions. *Journal of the ACM*  
483 (*JACM*), 21(2):201–206.  
484  
485 Chen, T., Xu, B., Zhang, C., and Guestrin, C. (2016). Training deep nets with sublinear memory  
cost. *arXiv preprint arXiv:1604.06174*.

---

486 Dao, T. (2023). Flashattention-2: Faster attention with better parallelism and work partitioning.  
487 *arXiv preprint arXiv:2307.08691*.  
488

489 Dao, T., Fu, D. Y., Ermon, S., Rudra, A., and Re, C. (2022). Flashattention: Fast and memory-  
490 efficient exact attention with IO-awareness. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho,  
491 K., editors, *Advances in Neural Information Processing Systems*.

492

493 Dean, J. and Ghemawat, S. (2008). Mapreduce: simplified data processing on large clusters. *Com-  
494 munications of the ACM*, 51(1):107–113.

495

496 Elliott, C. M. (2009). Beautiful differentiation. *ACM Sigplan Notices*, 44(9):191–202.

497

498 Gomez, A. N., Ren, M., Urtasun, R., and Grosse, R. B. (2017). The reversible residual network:  
499 Backpropagation without storing activations. *Advances in neural information processing systems*,  
500 30.

501 Griewank, A. and Walther, A. (2000). Algorithm 799: revolve: an implementation of checkpoint-  
502 ing for the reverse or adjoint mode of computational differentiation. *ACM Trans. Math. Softw.*,  
503 26(1):19–45.

504

505 Hsu, P.-L., Dai, Y., Kothapalli, V., Song, Q., Tang, S., Zhu, S., Shimizu, S., Sahni, S., Ning, H.,  
506 and Chen, Y. (2024). Liger kernel: Efficient triton kernels for llm training. *arXiv preprint  
507 arXiv:2410.10989*.

508

509 Rush, K., Charles, Z., and Garrett, Z. (2024a). Federated automatic differentiation. *Journal of  
510 Machine Learning Research*, 25(357):1–39.

511

512 Rush, K., Charles, Z., Garrett, Z., Augenstein, S., and Mitchell, N. (2024b). Drjax: Scalable and  
513 differentiable mapreduce primitives in jax. *arXiv preprint arXiv:2403.07128*.

514

515 Shoeybi, M., Patwary, M., Puri, R., LeGresley, P., Casper, J., and Catanzaro, B. (2019). Megatron-  
516 lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint  
517 arXiv:1909.08053*.

518

519 Wijmans, E., Huval, B., Hertzberg, A., Koltun, V., and Krähenbühl, P. (2024). Cut your losses in  
520 large-vocabulary language models. *arXiv preprint arXiv:2411.09009*.

## 521 A NOTATION

522

523 **Types & Terms**  $a : A$  denotes that the term  $a$  has type  $A$ . For example,  $a : \mathbb{R}$  denotes that the term  
524  $a$  is a real number.

525

526 **Functions**  $A \rightarrow B$  denotes the type of functions from  $A$  to  $B$ . Lambda functions are written out as  
527  $f = x \mapsto x^2 + 3$ , meaning  $f(x) = x^2 + 3$ .  
528  $\rightarrow$  is interpreted in a right-associative manner, meaning  $A \rightarrow B \rightarrow C = A \rightarrow (B \rightarrow C)$ ,  
529 which is isomorphic to  $A \times B \rightarrow C$ .

530 **Product Types**  $A \times B$  denotes the product type of  $A$  and  $B$ . For brevity, terms of product types  
531 are occasionally written out as tuples. If  $a : A$  and  $b : B$ , we can write  $(a, b) : A \times B$ ,  
532 and vice versa. For a term  $x : A \times B$ ,  $x_0$  is the term corresponding to the left-hand  
533 type ( $A$ ), and  $x_1$  corresponds to the right-hand type. We also used named product types,  
534 where  $\{a : A, b : B\}$  is the type of pairs of  $A$  and  $B$  indexed by the names  $a$  and  $b$ , i.e.  
535  $X : \{a : A, b : B\}$  says that  $X$  is a named product type, with two named parts  $a$  and  $b$   
536 of type  $A$  and  $B$  respectively. Similarly, terms of named product types are written out as  
537  $x = \{a : v, b : w\}$ , meaning  $x_a = v, x_b = w$ .

## 538 B BATCHED MONOIDREDUCE

---

540 Listing 2: A working, generalized, batched, Proof of Concept of MonoidReduce, implemented as a  
 541 torch.autograd.Function-factory.  
 542

```

543 import torch
544 from torch.autograd import Function
545 from torch.autograd.function import once_differentiable

546 def MonoidReduce(
547     fun_name ,
548     *,
549     init , # construct identity-tensors based on a and b
550     avals , # number of A-tensors
551     bvals , # number of B-tensors
552     islice , # slicing over i-axis
553     jslice , # slicing over j-axis
554     proj_fold , # fused projection and fold over A and B-slices .
555     proj_fold_bwd , # fused projection and fold backwards op.
556     binary_reduce , # monoidal product
557     ):
558     """
559     Constructor for MonoidReducers.
560     """

561     #aslice = slice(0, avals)
562     astop = avals
563     bstop = astop + bvals

564     class DynamicFunction(torch.autograd.Function):
565         @staticmethod
566         def forward(*inputs):
567             a = inputs[0:astop]
568             b = inputs[astop:bstop]
569             p = init(a, b)
570             # Running over ai can be done in parallel
571             # (requires rework of how to construct final p)
572             for ai, pi in zip(
573                 islice(a),
574                 islice(p),
575             ):
576                 pi_acc = pi
577                 # Each bj slice can be run in parallel, requires
578                 # parallel reduction of the pij-values over j.
579                 for bj in jslice(b):
580                     pij = proj_fold(ai, bj)
581                     pi_acc = binary_reduce(pi_acc, pij)

582             for pi_p, pi_acc_p in zip(pi, pi_acc):
583                 pi_p.copy_(pi_acc_p)

584             return p

585         @staticmethod
586         def setup_context(ctx, inputs, outputs):
587             ctx.save_for_backward(
588                 *inputs, *outputs
589             )

590         @staticmethod
591         @once_differentiable
592         def backward(ctx, *gp):
593             saved = ctx.saved_tensors

594             a = saved[:astop]
595             b = saved[astop:bstop]

```

---

```

594         p = saved[bstop:]
595
596         ga = [a_p.new_zeros(a_p.shape) for a_p in a]
597         amask = [a_p.requires_grad for a_p in a]
598         gb = [b_p.new_zeros(b_p.shape) for b_p in b]
599         bmask = [b_p.requires_grad for b_p in b]
600         # Running over ai can be done in parallel
601         # (requires atomic adds or other parallel sum over gaij)
602         for ai, gai, pi, gpi in zip(
603             islice(a),
604             islice(ga),
605             islice(p),
606             islice(gp),
607         ):
608             # Running over bi can be done in parallel
609             # (requires atomic adds or other parallel sum over gbij)
610             for bj, gbj in zip(
611                 jslice(b),
612                 jslice(gb),
613             ):
614                 lgaij, lgbij = proj_fold_bwd(ai, bj, pi, gpi)
615                 for gai_p, lgaij_p, requires_grad in zip(gai, lgaij, amask):
616                     if requires_grad:
617                         gai_p.add_(-lgaij_p)
618                 for gbj_p, lgbij_p, requires_grad in zip(gbj, lgbij, bmask):
619                     if requires_grad:
620                         gbj_p.add_(-lgbij_p)
621
622         return *ga, *gb
623
624     DynamicFunction.__name__ = f'MonoidReduce_{fun_name}'
625     DynamicFunction.__qualname__ = f'MonoidReduce_{fun_name}'
626     DynamicFunction.__module__ = getattr(init, '__module__', __name__)
627
628     return DynamicFunction

```

Listing 3: An implementation of attention using the MonoidReduce-factory from listing 2.

```

626 def init(a, b):
627     aq, = a
628     _, bv = b
629     B = aq.shape[0]
630     D = bv.shape[1]
631
632     z = aq.new_full((B,), float('-inf'))
633     v = bv.new_zeros((B, D))
634     return z, v
635
636 def islice(parts):
637     return zip(*(part.chunk(8) for part in parts))
638
639 def proj_fold(a, b):
640     aq, = a
641     bk, bv = b
642     logits = aq @ bk.t()
643     z = torch.logsumexp(logits, dim=1)
644     v = (logits - z[:, None]).exp() @ bv
645     return z, v
646
647 def binary_reduce(x, y):
648     xz, xv = x
649     yz, yv = y
650     z = torch.logaddexp(xz, yz)
651     v = (xz - z).exp()[:, None] * xv + (yz - z).exp()[:, None] * yv
652     return z, v

```

---

```

648
649 def proj_fold_bwd(a, b, p, gp):
650     # ignore gz, as it is only used for aggregation.
651     aq, bk, bv, pz, pv, _, gv = *a, *b, *p, *gp
652     # Recompute relevant stuff
653     logits = aq @ bk.t()
654     ws = (logits - pz[:, None]).exp()
655     # calculate gradients.
656     # (gv * pv).sum(1) could be precomputed before broadcast,
657     gbv = ws.t() @ gv
658     glogits = (gv @ bv.t() - (gv * pv).sum(1)[:, None]) * ws
659     gaq = glogits @ bk
660     gbk = glogits.t() @ aq
661
662     return (gaq,), (gbk, gbv)
663
664 Attention = MonoidReduce('Attention',
665                           init=init,
666                           avals=1,
667                           bvals=2,
668                           islice=islice,
669                           jslice=islice,
670                           proj_fold=proj_fold,
671                           proj_fold_bwd=proj_fold_bwd,
672                           binary_reduce=
673                           binary_reduce
674                           )
675
676 def attention(q, k, v):
677     """
678     q: L x DI
679     k: R x DI
680     v: R x DO
681     """
682     z, v = Attention.apply(q, k, v)
683     return v
684
685 Listing 4: An implementation of cross-entropy between two parameterized distributions using the
686 MonoidReduce-factory from listing 2.
687
688 def init(a, b):
689     ap, _ = a
690     B = ap.shape[0]
691     p = ap.new_full((B,), float('-inf'))
692     q = ap.new_full((B,), float('-inf'))
693     n = ap.new_zeros((B,))
694     return p, q, n
695
696 def islice(parts):
697     return zip(*(part.chunk(8) for part in parts))
698
699 def proj_fold(a, b):
700     ap, aq = a
701     bc, bd = b
702     pl = ap @ bc.t()
703     ql = aq @ bd.t()
704
705     p = torch.logsumexp(pl, dim=1)
706     q = torch.logsumexp(ql, dim=1)
707     n = ((ql - q[:, None]).exp() * pl).sum(1)
708     return p, q, n
709
710 def binary_reduce(x, y):
711     xp, xq, xn = x
712     yp, yq, yn = y

```

```

702     p = torch.logaddexp(xp, yp)
703     q = torch.logaddexp(xq, yq)
704     n = xn * (xq - q).exp() + yn * (yq - q).exp()
705     return p, q, n
706
707 def proj_fold_bwd(a, b, p, gp):
708     (lp, lq, ln), local_vjp = vjp(proj_fold, a, b)
709
710     p, q, n = p
711     gp, gq, gn = gp
712     # Compute  $d p^{ij} / d [a, b]$ .
713
714     lgp = gp * (lp - p).exp()
715     lgq = (gq + gn * (ln - n)) * (lq - q).exp()
716     lgn = gn * (lq - q).exp()
717     lg = (lgp, lgq, lgn)
718
719     # Run local vjp.
720     return local_vjp(lg)
721
722 XEntropy2 = MonoidReduce(
723     'XEntropy2',
724     init=init,
725     avals=2,
726     bvals=2,
727     islice=islice,
728     jslice=islice,
729     proj_fold=proj_fold,
730     proj_fold_bwd=proj_fold_bwd,
731     binary_reduce=binary_reduce)
732
733 def xentropy2(p, c, q, d):
734     p, _, n = XEntropy2.apply(p, q, c, d)
735     return p - n

```

The code presented in Listing 2, while functional, should be regarded as a Proof Of Concept. It is neither parallel nor very efficient. It does, however, showcase the memory-efficiency aspects: Figures 3 and 4 show memory profiles of attention (Listing 3) using the above implementation (extracted using `torch.profiler` with `cpu-pytorch`).

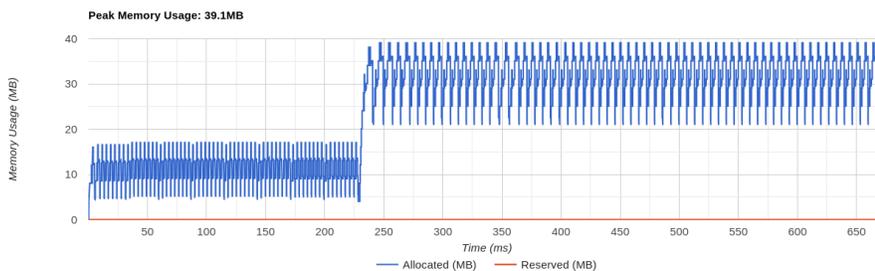


Figure 3: Memory profile for MonoidReduce attention (forward + backward).

## C MONOIDS CORRESPONDING TO CROSS ENTROPY AND TWO-LAYER MLPs

751  
752  
753  
754  
755

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

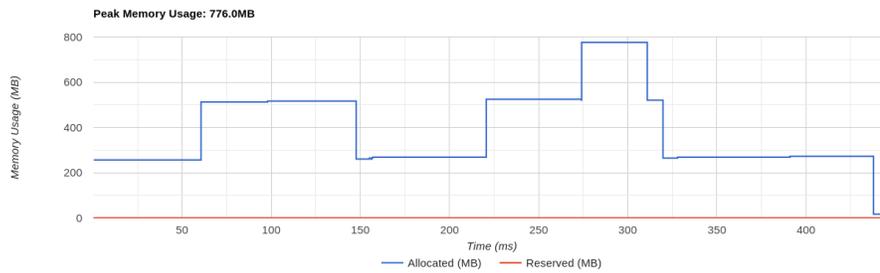


Figure 4: Memory profile for naive attention:  $\text{softmax}(q @ k.t()) @ v$ , (forward + backward)

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Input	Map	Reduce
<u>Attention: <math>Y = \text{softmax}(QK^T)V</math></u>		
$Q : \mathbb{R}^{M \times F}$	$T : \{z : \mathbb{R}, v : \mathbb{R}^D\}$	$P : T^M$
$K : \mathbb{R}^{N \times F}$	$A_i = Q_i$	$P_i = \bigodot_j f(A_i, B_j)$
$V : \mathbb{R}^{N \times D}$	$B_j = \{k : K_j, v : V_j\}$	
$\downarrow$	$f(a, b) = \left\{ \begin{array}{l} z : \langle a, b_q \rangle \\ v : b_v \end{array} \right\}$	$a \odot b = \left\{ \begin{array}{l} z : \ln(e^{az} + e^{bz}) \\ v : a_v e^{a_z - z} + b_v e^{b_z - z} \end{array} \right\}$
$Y : \mathbb{R}^{M \times D}$		$Y_i = P_{iv}$
<u>Cross Entropy: <math>Y = \text{cross-entropy}(\text{logits} = PC^T, \text{targets} = T)</math></u>		
$P : \mathbb{R}^{M \times D}$	$T : \{p : \mathbb{R}, n : \mathbb{R}\}$	$P : T^M$
$C : \mathbb{R}^{N \times D}$	$A_i = \{p : P_i, t : T_i\}$	$P_i = \bigodot_j f(A_i, B_j)$
$T : N^M$	$B_j = \{c : C_j, ix : j\}$	
$\downarrow$	$f(a, b) = \left\{ \begin{array}{l} p : \langle a_p, b_c \rangle \\ n : a_t = b_{ix} ? p : 0 \end{array} \right\}$	$a \odot b = \left\{ \begin{array}{l} p : \ln(e^{a_p} + e^{b_p}) \\ n : a_n + b_n \end{array} \right\}$
$Y : \mathbb{R}^M$		$Y_i = P_{ip} - P_{in}$
<u>Cross Entropy: <math>Y = \text{cross-entropy}(\text{logits} = PC^T, \text{targets} = \text{softmax}(QD^T))</math></u>		
$P : \mathbb{R}^{M \times D}$	$T : \{p : \mathbb{R}, q : \mathbb{R}, n : \mathbb{R}\}$	$P : T^M$
$C : \mathbb{R}^{N \times D}$	$A_i = \{p : P_i, q : Q_i\}$	$P_i = \bigodot_j f(A_i, B_j)$
$Q : \mathbb{R}^{M \times E}$	$B_j = \{p : C_j, q : D_j\}$	
$D : \mathbb{R}^{N \times E}$	$f(a, b) = \left\{ \begin{array}{l} p : \langle a_p, b_p \rangle \\ q : \langle a_q, b_q \rangle \\ n : p \end{array} \right\}$	$a \odot b = \left\{ \begin{array}{l} p : \ln(e^{a_p} + e^{b_p}) \\ q : \ln(e^{a_q} + e^{b_q}) \\ n : a_n e^{a_q - q} + b_n e^{b_q - q} \end{array} \right\}$
$\downarrow$		
$Y : \mathbb{R}^M$		$Y_i = P_{ip} - P_{in}$
<u>MLP: <math>Y = \sigma(XP^T)Q</math></u>		
$X : \mathbb{R}^{B \times M}$	$T : \mathbb{R}^N$	$P : T^B$
$P : \mathbb{R}^{K \times M}$	$A_i = X_i$	$P_i = \bigodot_j f(A_i, B_j)$
$Q : \mathbb{R}^{K \times N}$	$B_j = \{p : P_j, q : Q_j\}$	
$\downarrow$	$f(a, b) = \sigma(\langle a, b_p \rangle) b_q$	$a \odot b = a + b$
$Y : \mathbb{R}^{B \times N}$		$Y_i = P_i$

Table 2: Examples of Monoids  $T$ , map functions from inputs to intermediate values, reduction operations, and out projections that realize different operations used in Artificial Neural Networks.

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

Monoid	Gradient ( $D$ )
<u>CrossEntropy</u>	
$T : \{p : \mathbb{R}, n : \mathbb{R}\}$ $a \odot b = \left\{ \begin{array}{l} p : \ln(e^{a^p} + e^{b^p}) \\ n : a_n + b_n \end{array} \right\}$ $= c$	$\frac{d c}{d a} = g \mapsto \left\{ \begin{array}{l} p : g_p e^{a^p - c^p} \\ n : g_n \end{array} \right\}$
<u>CrossEntropy between two distributions</u>	
$T : \{q : \mathbb{R}, p : \mathbb{R}, n : \mathbb{R}\}$ $a \odot b = \left\{ \begin{array}{l} p : \ln(e^{a^p} + e^{b^p}) \\ q : \ln(e^{a^q} + e^{b^q}) \\ n : a_n e^{a^q - q} + b_n e^{b^q - q} \end{array} \right\}$ $= c$	$\frac{d c}{d a} = g \mapsto \left\{ \begin{array}{l} p : g_p e^{a^p - c^p} \\ q : (g_q + g_n (a_n - c_n)) e^{a^q - c^q} \\ n : g_n e^{a^q - c^q} \end{array} \right\}$
<u>MLP</u>	
$T : \{v : \mathbb{R}^D\}$ $a \odot b = \{v : a_v + b_v\}$ $= c$	$\frac{d c}{d a} = g \mapsto g$

Table 3: Proofs that condition 1 holds for the given monoids, and by extension that Theorem 2.1 and all its corollaries hold.