
Dali: Scaling Lazy Compilation & JIT Fusion

Jonathan Raiman
Dali
San Francisco, CA
jonathan@dali.ml

Abstract

The growing computational needs of many large scale machine learning applications have motivated the need for better JIT Fusion and graph manipulation tools that are able to improve utilization or stretch the on-chip memory capabilities through algorithmic tricks such as pipelining and recomputation. While there is strong evidence that these techniques can help keep up with the computational demands, they are typically hand written and require domain expertise in the internals of machine learning frameworks, GPU programming, and computational graph manipulation. Recent work, such as AutoTVM, FlexFlow, or Dali have shown that many of these optimizations can be automatically discovered and applied by the framework without human labor or expertise. Unfortunately, many of the proposed algorithms cannot handle fusion across entire deep networks or require hours of auto-tuning, and as we remark in this paper, choosing what operations to fuse is equally challenging. We propose to improve Dali’s compiler through the addition of a more robust fusion boundary detector and variable elimination during code generation. We demonstrate the impact of these changes by showing how JIT Fusion and graph optimizations are able to accelerate deep Transformers and RNNs. We compare our approach to TensorFlow, TensorFlow with XLA, and PyTorch, and observe a $1.23 - 2.13\times$ speedup over TensorFlow, and $1.08 - 1.92\times$ speedup over PyTorch.

1 Introduction

The exponential rise in computational requirements in the largest machine learning experiments [1] presents an important challenge for future machine learning systems that wish to keep up with demand. One solution to improving hardware utilization comes through the use of optimizations such as pipelining [2] or gradient checkpointing [3, 4], and the use of task-optimized code via autotuning [5] or model-based code generation [6, 7, 8]. In order to benefit from these improvements, human labor and expertise is needed to implement pipelining, tradeoff memory and computation when recomputing, while autotuning requires hours of profiling [7, 5, 9]. Learning or specifying a model has emerged as a portable and fast option for adapting optimizations to new cases without human intervention or significant retraining or search time [10, 7, 6]. However, the search-based techniques in Dali [6] have only been demonstrated on static CNNs, where a greedy JIT Fusion strategy can lead to poor operation coupling and slower execution. Furthermore, code generation relies on A* [11] which has an exponential complexity relative to the depth of the search.

In this work, we propose to improve Dali’s optimization capabilities by introducing a smarter JIT Fusion strategy. We also enable Dali to discover and eliminate redundant variables during code generation to shrink the search space. We empirically validate the scalability of these changes on a CNN, Transformer, and mLSTM. Concretely our key contributions are the following:

1. We show how to scale JIT Fusion and code generation to larger models and find that this approach now outperforms PyTorch and TensorFlow. A new fusion strategy lets us find useful boundaries in the computation graphs. This change enables support for deeper networks while finding practical boundaries to prevent naive coupling of operations.
2. We add discovery and elimination of loop variables to code generation A* [11] search process in Dali, enabling it to shrink the solution search space by several orders of magnitude and handle the larger kernels found in Transformers and mLSTMs.
3. We provide results showing cached graph optimization reduces by a factor of 20 – 257× the transformation cost, suggesting graph optimizations are applicable to dynamic graphs ranging from static CNNs to Transformers and RNNs.

2 A* Fused CUDA Kernel Generation

2.1 Approach overview



Figure 1: Adjacent operations with the ability to be JIT compiled can be joined into a single compilation unit via Fusion.

Fusing operations together and generating task-specific code is a powerful tool for accelerating computation graphs by removing temporary variables and reducing slow memory movement [12, 13, 14]. In this work we build upon the JIT Fusion optimizations in Dali’s compiler [6], which rely on having each operation indicate whether they are capable of generating code that can be concatenated and compiled into a larger fused kernel. In this prior work, JIT Fusion operates in two stages: (1) combine a series of operations into a larger compilation unit, (2) make parallelism and design decisions during code generation to compile a fast kernel.

In the first stage, to find a JIT subgraph to compile, Dali considers connected subtrees of the full computation graph composed of JIT-able expressions as depicted in Figure 1. If the parent of a JIT-able expression is also JIT-able it will combine the two operations into a single computation unit through *fusion*. The combined set of operations can now be thought of as a single operation. This process is repeated until all JIT-able expressions are fused to their neighbors.

The second stage collects all loops and code design choices from the nodes inside the JIT subgraph. A model of the GPU’s memory and parallelism informs an A* search [11] through different ways of designing the code ranked by a cost function similar to the approach taken by [15, 16].

Through this approach it is possible without human intervention to obtain efficient code for operations composed of many primitives such as Softmax, Layer Norm [17], Batch Norm [18], Attention [19, 20], etc. However, as we remark in this work (Figure 2), combining operations in a greedy fashion (stage 1) leads to two significant problems: (1) the code generation (stage 2) solution space grows exponentially with the number of variables, and (2) the generated code becomes over-constrained and cannot be parallelized.

We address these issues in Section 2.2 where we introduce a Fusion strategy that decides which nodes to fuse, and in Section 2.3 we describe our approach to discovering and eliminating redundant variables during code generation to scale to larger kernels.

2.2 Fusion strategy

When constructing computation graphs, all parent-child relations involving JIT-able operations are candidates for *fusion*. Using *fusion*, we can write Softmax, Attention, CrossEntropy, or any other

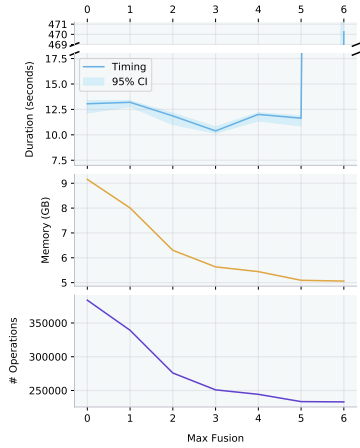


Figure 2: Impact of F_{\max} heuristic on runtime and memory when training a Transformer.

Table 1: Runtime & Memory Usage

Framework	Time (s) / Epoch ($\mu \pm \sigma$)	MB / Step
Task		
CNN - MNIST		
Dali w/o JIT Fusion	1.29 ± 0.0050	376.62
Dali w/o temp	1.13 ± 0.0046	363.38
Dali	1.11 ± 0.0034	363.38
PyTorch	2.14 ± 0.0037	913.00
TensorFlow	2.37 ± 0.0239	1178.00
TensorFlow + XLA	1.53 ± 0.0368	959.00
Task		
Transformer - 1 Billion Word		
Dali w/o JIT Fusion	12.518 ± 0.5702	9165.89
Dali	8.735 ± 0.0464	4276.50
PyTorch	9.476 ± 0.2425	4003.76
TensorFlow	11.123 ± 0.1893	4059.63
Task		
mLSTM - Amazon Reviews		
Dali w/o JIT Fusion	95.326 ± 0.1528	9581.12
Dali w/o temp	268.382 ± 1.5226	7094.12
Dali	93.342 ± 0.0889	7102.13
PyTorch	114.868 ± 0.0240	4967.00
TensorFlow	115.120 ± 0.1336	4363.00

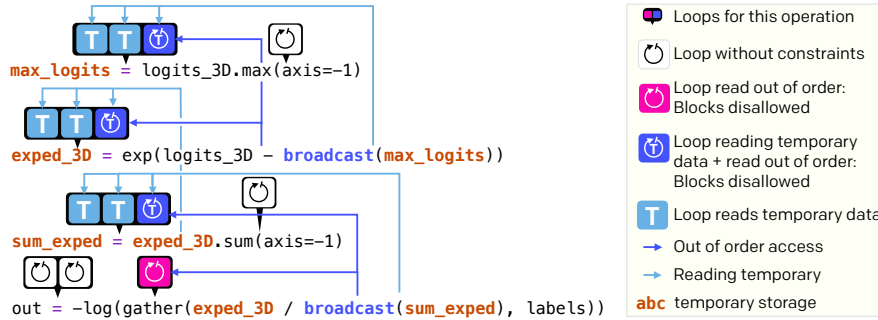


Figure 3: Constraint graph for `sparse_softmax_cross_entropy(logits, labels)` with 3D logits. Broadcast operations perform out of order access on their input: `max_logits` and `sum_exped`. Because a different block might have written to these inputs during the evaluation of the kernel, we disallow using “block parallelism” in the loops connected to broadcast operations.

numerical operations using primitives such as addition, subtraction, division, reductions, etc. and let the code generation system deal with how to plan, reuse, store, and combine these operations to compute the result. However, when planning our computation, we discover many constraints on the parallelism options for each reduction or assignment loop. In cases where we choose to fuse many operations together, we may end up in a situation with constraints that force the creation of a very inefficient function.

In many cases, logical operation groups such as loss computation, normalization, or affine transformations involve 1-3 inputs. However, an entire computation graph might actually involve 100s of unique weight matrices, many temporary variables, and many other noise inputs when using functions such as Dropout [21]¹. Because many of the operations in a computation graph might be JIT-able, it is possible to attempt to fuse 10s to 100s of operations into a single compilation unit. This has several downsides: (1) with more stages in the computation, new constraints on parallelism arise, leading in some cases to non-performant code, (2) the additional loops and variables increase the search space for code generation dramatically, a problem we discuss further in Section 2.3, (3) there is a maximum number of data and arguments that can be given to a CUDA kernel, placing an upper

¹We did not explore fusing random number generators, but it is possible (e.g. Tausworthe Generators [22]).

bound on the largest fused kernel, or requiring a leaner kernel call strategy. For these reasons it is crucial to discover good fusion boundaries. We propose a simple heuristic that can be computed greedily:

1. Fuse JIT operations bottom-up; record how many unique input arrays they require.
2. If the fusion of a new JIT operation would increase the number of unique inputs past the maximum number F_{\max} , create a new compilation unit with this operation as a leaf. Otherwise, continue extending.

We measure empirically the effect of F_{\max} on runtime, memory consumption, and number of executed operations when training a Transformer [20] and report our findings in Figure 2. As expected, increasing F_{\max} leads to more aggregation of JIT operations, and a decrease in the number of executed operations. Fusion also decreases memory consumption by allowing chaining of operations without temporary variables. As described earlier, fusing too many operations together can lead to non-performant code as is visible with $F_{\max} = 6$, where runtime jumps from 11 to 470 seconds. Surprisingly, we observe a Goldilocks principle in the runtime: setting F_{\max} to 3 gives the best results².

2.3 Variable Coupling and Elimination

A common problem Dali faces when handling large scale models is generating kernels involving many loops, since loop count exponentially grows the search space. Upon inspection, many of the variables in these loops have constraints and dependencies that make them redundant, thus by coupling their assignments with other variables we can greatly reduce the search space.

We see this problem arise when generating code for an 8-layer Transformer [20] if we do not set F_{\max} and allow all neighboring JIT operations to combine. During code generation we encounter a single kernel involving 36 loops and 12 assignments (search space size: 8,796,093,022,208). However, many of the variables in this example are implicitly coupled to other variables: e.g. two iterations over the same batch dimension of an array will be constrained to be data-consistent (so the first loop has to operate on threads or the same block as the second loop). The most efficient use of our GPU resources will involve these loops sharing the same parallelism strategy (e.g. use blocks or threads in both cases). Thus, we can look at all parent-child relationships in our constraint graph and derive which variables are actually redundant³. With this information, Dali eliminates 17 variables and shrinks the search space $8 \cdot 10^6$ times down to 1,048,576, making it possible to generate the kernel.

2.4 Constraint Extraction Example: Softmax Cross Entropy

To illustrate the code generation strategy, let us take a commonly used operation: computing the cross entropy between a set of labels and a generated unnormalized probability distribution. To compute this, we commonly chain `softmax(logits)` and `sparse_cross_entropy(probs, labels)`. Concretely, using Dali:

```
Array sparse_softmax_cross_entropy(const Array& logits, const Array& labels) {
    auto exped = (logits - logits.max({-1}, /*keepdims=*/true)).exp();
    auto probs = exped / exped.sum({-1}, /*keepdims=*/true); // softmax
    return -op::log(op::gather_from_rows(probs, labels)); /* cross-entropy*/
}
```

Only two inputs are needed, satisfying the F_{\max} condition for fusion boundaries. As a result, Dali combines all the operations into one. After observing broadcast and reuse, Dali stores the result of `max`, `exp`, and `sum` into temporary variables. The final computation has four assignments: three temporary variables and one output. The assignments and reductions involve a total of 14 loops (2 reductions, 1 loop controlling `gather_from_rows` iteration, and 11 used for assignment). Dali annotates these loops with data dependencies that induce constraints on the generated code as illustrated in Figure 3.

²We see this effect repeated for the mLSTM and CNN.

³Redundancy detection checks if two loops are descendants of each other, have the same symbolic size, and whether they have equal access to using threads or blocks.

Table 2: Optimization Time

Step	CNN	Model	mLSTM
		Transformer	
Percent Total			
Optimization	0.64%	0.11%	0.052%
Cached Transformation	1.88%	1.88%	2.69%
Expression hash	2.33%	1.56%	1.71%
Computation	95.85%	96.44%	95.54%
Time/Call ($\mu \pm \sigma$)			
Optimization	41.11ms \pm 320 μ s	183.3ms \pm 1.91ms	794.9ms \pm 10.3ms
Cached Transformation	106.31 μ s \pm 1.00 μ s	3.041ms \pm 8.50 μ s	40.56ms \pm 150 μ s
Expression hash	53.36 μ s \pm 0.160 μ s	2.520ms \pm 7.03 μ s	25.80ms \pm 74.5 μ s
Computation	4.366ms \pm 14.3 μ s	155.8ms \pm 3.25ms	1.439s \pm 57.8ms

The search space has 4,194,304 possible solutions, but by observing that certain loops can be coupled, Dali reduces the number of variables to assign from 14 to 7, and the search space drops by **8,192 times** down to 512. Dali searches through the space of solutions and evaluate 155 candidates, and arrives at a solution that uses blocks for the leading two “batch” dimensions, and threads everywhere else. This solution matches the one written by human programmers in PyTorch and TensorFlow.

3 Results

We investigate the scalability and benefits of using a fusion strategy and variable elimination through a C++ implementation of the Dali compiler [6] on a series of neural network training tasks. We summarize our timing experiments regarding graph transformation on 3 tasks and architectures: image classification with a CNN on MNIST (28x28 grayscale), subword language modeling with a Transformer [20] on the 1 Billion Word dataset [23]⁴, and character language modeling using the mLSTM from [25] on the Amazon Reviews dataset introduced by [26]. We measure memory usage and runtime per epoch over (1) 100 epochs of training a CNN, (2) 10 epochs of training a Transformer, and (3) 10 epochs of an mLSTM. Full hyperparameters can be found in Appendix A.

All experiments are run on a 12-core 3.60GHz Intel i7-6850K CPU with an NVIDIA Titan X Pascal. We compare six configurations: Dali with and without JIT Fusion, Dali with no temporary storage in the generated code, TensorFlow 1.13.1 with and without XLA⁵, and PyTorch 1.1 [27] with all frameworks using CUDA 10.1 and CuDNN 7.5.1 primitives [28].

We report our results in Table 1 and find that we always obtain a speedup over TensorFlow and PyTorch. Specifically we see a 1.92/2.13 \times speedup over PyTorch/TensorFlow with a CNN, a 1.23/1.23 \times speedup over PyTorch/TensorFlow with an mLSTM, and a 1.08/1.27 speedup over PyTorch/TensorFlow using a Transformer. We also observe that using temporary storage in the generated code requires a small amount of memory, but can significantly speedup runtime.

We also time fine-grained steps of the graph optimization in Table 2. In those measurements, we note that hashing and cached transformation is \approx 20-257 \times faster than the initial optimization pass. The optimization overhead remains modest across static models (CNN) and dynamic and larger models (Transformer, mLSTM).

4 Related Work

Our approach to optimizing computation graphs specifically for machine learning applications is closely related to work on powerful DSLs and code-generation tools that generate, compose, or design high performance code. These can generally be divided into two groups: (1) high-level specification of specific operations such as convolutions used to automatically generate kernels, (2) functions in

⁴We use the `language_model_1m1b32k` problem from `tensor2tensor` [24].

⁵XLA support is currently experimental, so we were only able to obtain results with XLA on the CNN example. The Transformer and mLSTM segfault after compilation.

scripting languages, program traces, or syntax trees being converted into portable and more efficient compiled code.

There is a vast amount of work in group (1) focusing on Polyhedral methods and integer linear programming to optimize cost [29, 30], or exposing different parallelism and tiling controls using a DSL such as Halide and TVM [15, 16] to generate high performance code for specific hardware. Tensor comprehensions [5] builds upon this work by converting high-level kernel specifications into Halide code and autotuning it, while in AutoTVM [7] the authors learn a cost estimator to guide code generation. In this work, we construct a similar representation as Halide/TVM, but extend this work by enabling constraint extraction across multiple stages of computation within a CUDA kernel. This enables higher expressivity, the use of temporary memory, and a greater ability to combine JIT-able operations into a single kernel. We differ from prior autotuning work in our use of an A* search algorithm to design the code: we prune out large portions of the search space by using a cost function with an admissible heuristic. We view learning based-approaches such as AutoTVM [7], automatic device placement [31, 32, 33], or learnt fusion strategies [34] as complementary to ours, because they can augment our compiler with smarter optimization passes.

Within group (2) we see a mix of tracing based approaches such as PyTorch scripting [35], syntax tree extraction such as MatchBox and AutoGraph [36, 37], and a combination of tracing and JIT compilation within JAX [14]. These methods have been successful at removing the overhead of scripting languages, creating portable modules, and in some cases converting control flow into efficient element-wise operations. A common goal of these approaches is to convert user imperative code into an internal computation graph that can be optimized and executed more efficiently, usually through the use of decorators or source code annotation. Our work resembles the tracing approaches, however we build a dynamic computation graph which is lazily compiled and executed without any extra annotation or decorators, similar to Dynet’s autobatching [38], and our optimized graph is dimension-agnostic, covers multiple optimization passes, and generates fused code most similar to the transformation passes found in JAX. In our approach, by ignoring function boundaries, we benefit by potentially discovering interesting optimizations within a larger program, which can later be reapplied to inputs with different dimensions. However, a downside of tracing is the inability to observe native control flow: Dali instead receives unrolled loops, thereby requiring new Transformation Graphs for similar input graphs.

5 Conclusion

We have shown how improvements to Dali’s compiler allow it to scale to deeper networks and outperform existing computational frameworks on a CNN, Transformer, and mLSTM. These gains can be attributed to two new capabilities for scaling JIT Fusion and graph optimizations to large computational graphs where the abundance of fusion opportunities can cripple greedy computation graph compilers. The first capability involves the addition of JIT Fusion boundaries based on a simple heuristic that accelerates the execution of a CNN, Transformer, and mLSTM. We hypothesize that controlling the number of input arrays in a kernel is correlated with the likelihood of introducing incompatible parallelism constraints, which would force the generation of slower code. The second capability is the discovery and elimination of variables in the kernel prior to performing an A* search for the most parallel solution. In prior work [6], the number of variables could grow to 50 or more, leading to intractably large search spaces. We show how detecting that loop variables are redundant can enable us to shrink the search space by 6 orders of magnitude, enabling code generation to be usable in deeper networks where more operations are fused.

As future work, we have shown the importance of fusion boundaries, and we expect a data-driven or learnt policy following [32, 31, 7] would perform even better at choosing which operations to fuse. While we have shown that there algorithmic approaches to further accelerate cached optimization transformation, an important challenge is knowing how to detect similar subgraphs to improve optimization reuse across computation graphs⁶. Lazy compilation makes computation of one-off large graphs impractical since the optimization time is never amortized by reuse. Many of the optimizations applied by Dali are local, thus it is conceivable that previous optimization passes could accelerate unseen graph evaluation.

⁶For instance we should be able to reuse optimization passes across computation graphs involving RNNs unrolled for a different number of timesteps by detecting the temporal loop.

References

- [1] Dario Amodei and Danny Hernandez. AI and Compute. <https://openai.com/blog/ai-and-compute/>, 2018.
- [2] Yanping Huang, Yonglong Cheng, Dehao Chen, HyoukJoong Lee, Jiquan Ngiam, Quoc V Le, and Zhifeng Chen. Gpipe: Efficient training of giant neural networks using pipeline parallelism. *arXiv preprint arXiv:1811.06965*, 2018.
- [3] Andreas Griewank and Andrea Walther. *Evaluating derivatives: principles and techniques of algorithmic differentiation*, volume 105. Siam, 2008.
- [4] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. *arXiv preprint arXiv:1604.06174*, 2016.
- [5] Nicolas Vasilache, Oleksandr Zinenko, Theodoros Theodoridis, Priya Goyal, Zachary DeVito, William S Moses, Sven Verdoolaege, Andrew Adams, and Albert Cohen. Tensor comprehensions: Framework-agnostic high-performance machine learning abstractions. *arXiv preprint arXiv:1802.04730*, 2018.
- [6] Jonathan Raiman. Dali: Lazy compilation of dynamic computation graphs. In *Workshop on Systems for Machine Learning and Open Source Software at NeurIPS 2018*, 2018.
- [7] Tianqi Chen, Lianmin Zheng, Eddie Yan, Ziheng Jiang, Thierry Moreau, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. Learning to optimize tensor programs. In *Advances in Neural Information Processing Systems*, pages 3389–3400, 2018.
- [8] Zhihao Jia, Sina Lin, Charles R Qi, and Alex Aiken. Exploring hidden dimensions in parallelizing convolutional neural networks. *arXiv preprint arXiv:1802.04924*, 2018.
- [9] Paul Barham and Michael Isard. Machine learning systems are stuck in a rut. In *Proceedings of the Workshop on Hot Topics in Operating Systems*, pages 177–183. ACM, 2019.
- [10] Ravichandra Addanki, Shaileshh Bojja Venkatakrishnan, Shreyan Gupta, Hongzi Mao, and Mohammad Alizadeh. Placeto: Efficient progressive device placement optimization. In *NIPS Machine Learning for Systems Workshop*, 2018.
- [11] Peter E Hart, Nils J Nilsson, and Bertram Raphael. A formal basis for the heuristic determination of minimum cost paths. *IEEE transactions on Systems Science and Cybernetics*, 4(2):100–107, 1968.
- [12] Haicheng Wu, Gregory Diamos, Srihari Cadambi, and Sudhakar Yalamanchili. Kernel weaver: Automatically fusing database primitives for efficient gpu computation. In *2012 45th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 107–118. IEEE, 2012.
- [13] Jiří Filipovič, Matúš Madzin, Jan Fousek, and Luděk Matyska. Optimizing cuda code by kernel fusion: application on blas. *The Journal of Supercomputing*, 71(10):3934–3957, 2015.
- [14] Roy Frostig, Matthew James Johnson, and Chris Leary. Compiling machine learning programs via high-level tracing, 2018.
- [15] Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman P. Amarasinghe. Halide: a language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines. In *PLDI*, 2013.
- [16] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Haichen Shen, Eddie Yan, Leyuan Wang, Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. Tvm: end-to-end compilation stack for deep learning. In *SysML Conference*, 2018.
- [17] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *stat*, 1050:21, 2016.
- [18] Sergey Ioffe. Batch renormalization: Towards reducing minibatch dependence in batch-normalized models. In *Advances in neural information processing systems*, pages 1945–1953, 2017.
- [19] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

- [21] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [22] Lee Howes and David Thomas. Efficient random number generation and application using cuda. *GPU gems*, 3:805–830, 2007.
- [23] Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. One billion word benchmark for measuring progress in statistical language modeling. *arXiv preprint arXiv:1312.3005*, 2013.
- [24] Ashish Vaswani, Samy Bengio, Eugene Brevdo, Francois Chollet, Aidan N. Gomez, Stephan Gouws, Llion Jones, Łukasz Kaiser, Nal Kalchbrenner, Niki Parmar, Ryan Sepassi, Noam Shazeer, and Jakob Uszkoreit. Tensor2tensor for neural machine translation. *CoRR*, abs/1803.07416, 2018.
- [25] Alec Radford, Rafal Jozefowicz, and Ilya Sutskever. Learning to generate reviews and discovering sentiment. *arXiv preprint arXiv:1704.01444*, 2017.
- [26] Julian McAuley, Rahul Pandey, and Jure Leskovec. Inferring networks of substitutable and complementary products. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 785–794. ACM, 2015.
- [27] Adam Paszke, Sam Gross, Soumith Chintala, and Gregory Chanan. Pytorch: Tensors and dynamic neural networks in python with strong gpu acceleration. *PyTorch: Tensors and dynamic neural networks in Python with strong GPU acceleration*, 6, 2017.
- [28] Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, Jonathan Cohen, John Tran, Bryan Catanzaro, and Evan Shelhamer. cudnn: Efficient primitives for deep learning. *arXiv preprint arXiv:1410.0759*, 2014.
- [29] Uday Bondhugula, Albert Hartono, Jagannathan Ramanujam, and Ponnuswamy Sadayappan. A practical automatic polyhedral parallelizer and locality optimizer. In *Acm Sigplan Notices*, volume 43, pages 101–113. ACM, 2008.
- [30] Sven Verdoolaege, Juan Carlos Juega, Albert Cohen, Jose Ignacio Gomez, Christian Tenllado, and Francky Catthoor. Polyhedral parallel code generation for cuda. *ACM Transactions on Architecture and Code Optimization (TACO)*, 9(4):54, 2013.
- [31] Azalia Mirhoseini, Hieu Pham, Quoc Le, Mohammad Norouzi, Samy Bengio, Benoit Steiner, Yuefeng Zhou, Naveen Kumar, Rasmus Larsen, and Jeff Dean. Device placement optimization with reinforcement learning. 2017.
- [32] Azalia Mirhoseini, Anna Goldie, Hieu Pham, Benoit Steiner, Quoc V Le, and Jeff Dean. A hierarchical model for device placement. 2018.
- [33] Ravichandra Addanki, Shaileshh Venkatakrishnan, Shreyan Gupta, Hongzi Mao, and Mohammad Alizadeh. Placeto: Efficient progressive device placement optimization. In *ML For Systems at NeurIPS 2018*, 2018.
- [34] Amirali Abdolrashidi, Qiumin Xu, Shibo Wang, Sudip Roy, and Yanqi Zhou. Learning to fuse. In *NeurIPS Machine Learning for Systems Workshop*, 2019.
- [35] Torch Contributors. Torch Script. <https://pytorch.org/docs/master/jit.html>, 2018. [Online; accessed 19-October-2018].
- [36] James Bradbury. Matchbox: Automatic Batching for Dynamic Deep Learning. <http://on-demand.gputechconf.com/gtc/2018/presentation/s8977-matchbox-automatic-batching-for-dynamic-deep-learning.pdf>, 2018. [Online; accessed 19-October-2018].
- [37] Tensorflow. AutoGraph: Easy control flow for graphs. <https://www.tensorflow.org/guide/autograph>, 2018. [Online; accessed 19-October-2018].
- [38] Graham Neubig, Yoav Goldberg, and Chris Dyer. On-the-fly operation batching in dynamic computation graphs. In *Advances in Neural Information Processing Systems*, pages 3971–3981, 2017.

A Hyperparameters

In this section we include the CNN architecture in Table 3, and hyperparameters for the experiments: CNN in Table 4, Transformer in Table 5, and mLSTM in Table 6.

Table 3: CNN Architecture

Layer	Window/Strides	Input channels	Output channels
Convolution + Relu	(5, 5)/(1, 1)	1	64
MaxPool	(2, 2)/(2, 2)	64	64
Conv + Relu	(5, 5)/(1, 1)	64	64
Convolution	(2, 2)/(2, 2)	64	64
FC + Relu		3136	1024
FC + Softmax		1024	10

Table 4: CNN Hyperparameters

Dataset	MNIST 28x28
Batch Size	256
Optimizer	SGD

Table 5: Transformer Hyperparameters

Dataset	1 Billion Word
Timesteps	100
Batch Size	32
Optimizer	SGD
Hidden Size	512
Intermediate Size	1024
Attention Heads	4
Number of Layers	8
Activation	GELU [1]
Examples/epoch	2048
Steps/epoch	64

Table 6: mLSTM Hyperparameters

Dataset	Amazon Reviews
Timesteps	256
Batch Size	32
Optimizer	SGD
Hidden Size	4096
RNN Weight Norm	Yes
Output Weight Norm	No
Vocabulary Size	256
Embedding Size	64
Examples/epoch	2048
Steps/epoch	64

References

- [Appendix1] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.