

MIOpen: An Open Source Library For Deep Learning Primitives

Jehandad Khan, Paul Fultz, Artem Tamazov, Daniel Lowell, Chao Liu,
Michael Melesse, Murali Nandhimandalam, Kamil Nasyrov, Ilya Perminov, Tejash Shah,
Vasilii Filippov, Jing Zhang, Jing Zhou, Bragadeesh Natarajan, Mayank Daga
AMD Inc.

Mayank.Daga@amd.com

Abstract—Deep Learning has established itself to be a common occurrence in the business lexicon. The unprecedented success of deep learning in recent years can be attributed to: abundance of data, availability of gargantuan compute capabilities offered by GPUs, and adoption of open-source philosophy by the researchers and industry. Deep neural networks can be decomposed into a series of different operators. MIOpen, AMD’s open-source deep learning primitives library for GPUs, provides highly optimized implementations of such operators, shielding researchers from internal implementation details and hence, accelerating the time to discovery. This paper introduces MIOpen and provides details about the internal workings of the library and supported features.

MIOpen innovates on several fronts, such as implementing fusion to optimize for memory bandwidth and GPU launch overheads, providing an auto-tuning infrastructure to overcome the large design space of problem configurations, and implementing different algorithms to optimize convolutions for different filter and input sizes. MIOpen is one of the first libraries to publicly support the `bfloat16` data-type for convolutions, allowing efficient training at lower precision without the loss of accuracy.

Index Terms—Convolution, Deep Learning, GPU, HIP, Machine Learning, MIOpen, OpenCL[®], Performance

I. INTRODUCTION

Deep Learning has burgeoned into one of the most important technological breakthroughs of the 21st century. The use of deep learning has garnered immense success in applications like image and speech recognition, recommendation systems, and language translation. This in turn advances fields like autonomous driving and disease diagnosis [1]. GPUs have played a critical role in the advancement of deep learning. The massively parallel computational power of GPUs has been influential in reducing the training time of complex deep learning models hence, accelerating the time to discovery [2]. The availability of open-source frameworks like TensorFlow and PyTorch is another cornerstone for the fast-paced innovation in deep learning [3], [4].

The deep learning frameworks decompose the models as either a computational graph or a sequence of operations [5]–[7]. These high-level operations are then compiled down to a series of hardware specific high-performance primitives. These primitives in deep learning are akin to BLAS (Basic Linear Algebra Subprograms) [8] in linear algebra and high performance computing. Availability of a library which provides highly optimized implementations of such primitives enables the deep learning researchers to focus on their science and leaves the

burden of developing such primitives on the hardware vendors. The library then provides a simple and callable application programming interface (API) to enable seamless integration with client libraries and be flexible so that new features may be added easily.

MIOpen is AMD’s deep learning primitives library which provides highly optimized, and hand-tuned implementations of different operators such as *convolution*, *batch normalization*, *pooling*, *softmax*, *activation* and layers for *Recurrent Neural Networks (RNNs)*, used in both training and inference [9]. Moreover, MIOpen is fully open-source including all its GPU kernels; complementing AMD’s open-source ROCm stack [10]. MIOpen is the first to extend the open-source advantage into GPU vendor libraries thereby, continuing to embark on the same ethos as the deep learning community.

As deep learning has gained critical acclaim over the years, substantial research has been conducted to accelerate it. One optimization technique called *fusion* has been recognized to be more potent than others [11]. Fusion allows to fuse or collapse several neural network layers thereby, optimizing on 1) memory bandwidth requirements by requiring less data to be moved between host and GPU memories, and 2) GPU kernel launch overheads by launching fewer GPU kernels compared to the vanilla, non-fused neural network. Aside from discrete primitives, MIOpen also offers a fusion API which allows the frameworks to fuse some of the operations mentioned above. MIOpen fusion can be used to accelerate both convolution and recurrent neural networks.

Another area that has flourished with the popularity of deep learning is open-source graph compilers [11], [12], [13], [14]. Graph compilers further the relevance of deep learning to wide-spread applications by generating the implementations of aforementioned operators instead of relying on hardware specific libraries. However, generating high-performance implementations of two operators, convolution and GEMM, is extremely cumbersome without inherent knowledge of the underlying hardware. Therefore, the graph compilers rely on libraries like MIOpen for these operators. MIOpen’s open-source nature enables a plethora of optimization opportunities which were not possible before. For example, fusing an operator generated by the compiler with MIOpen’s convolutions. MIOpen facilitates these optimization by breaking down complex operators like convolutions into several simple and small

operators and providing high-performance implementations of these simple operators to the graph compiler. This MIOpen feature is called *composable kernels*.

The primary aim of MIOpen is to provide access to high-performance kernels, support several data-types, and also support as many hardware targets as required. To that end, MIOpen supports four different data-types: `float32`, `float16`, `bfloat16`, and `int8`, and two programming models: HIP and OpenCL[®] [15]. The kernels in MIOpen are backed by both high-level language as well as hand-tuned assembly implementations. MIOpen also provides an auto-tuning infrastructure to achieve maximum performance on the user's hardware and software environment.

This document provides an under-the-hood look at the MIOpen library providing detailed information about the functionality of the library as well as introduce MIOpen's capabilities to users and developers. The rest of the paper is organized as follows: Section II describes some prior work, Section III describes the overall design philosophy of the library and provides details about kernel compilation, abstractions used to localize those details in the library, tuning infrastructure for improving kernel performance and MIOpen's support for OpenCL[®] and HIP. This is followed by Section IV which provides details about the supported operations; primarily the convolution operation. Section V describes MIOpen's Fusion API for merging different operations for increased performance, this is followed by some usage statistics and performance comparisons in Section VI. Section VII presents conclusion and future work.

II. RELATED WORK

Developing hardware-optimized libraries for most critical and time-sensitive operations is a well-known practice. For linear algebra such libraries are known as BLAS (Basic Linear Algebra Subsystem) and have different implementations for different systems [8], [16]–[18]. In similar spirit different deep learning libraries have been written, to make it easier for client applications to implement different deep learning primitives. Alex Kruschevsky's *cuda-convnet* is one of the initial libraries to implement convolutions and inspired many others [19], [20]. Chetlur et al. developed cuDNN, a deep neural network library for nVIDIA GPUs [21]. MIOpen falls in this category since it provides a C programming language based API for deep learning primitives. While these libraries aim to accelerate deep learning primitives on GPUs, research also been conducted to improve the performance of inference only loads on different CPUs such as MKL-DNN [20].

Most of the above mentioned libraries focus on lower level optimization opportunities. An orthogonal approach is to abstract this detail behind a domain specific language (DSL). This technique has already been successfully applied to other domains such as computer vision and linear algebra [22]–[25]. Vasilache et al. developed *Tensor Comprehensions*, which takes a similar approach and designs a language which can infer tensor dimensions and summation indices automatically [19]. However, such an approach makes it complicated

to support a wide array of platforms and hardware targets as is required of MIOpen.

A. MIOpen and higher level frameworks

The above libraries are augmented by a community of frameworks which enable researchers and practitioners to express their computation pipeline using a host language (typically Python[™] or some other higher level language) [5] [3] [26] [27] [28]. These frameworks in turn call out libraries such as MIOpen for efficient implementation of the primitives required to implement the computation in those graphs. Frameworks strive to support a wide array of hardware and applications, for instance both TensorFlow and PyTorch already support MIOpen as a backend aimed at AMD GPUs. Thus a user can seamlessly change the hardware target without changing their application code.

III. OVERALL DESIGN

This section describes the MIOpen's design philosophy using the convolution operation as an example.

A. Kernels and Solvers

Mapping a problem description to a particular kernel requires MIOpen to determine the file which contains the required GPU kernel, the name of the kernel in the file and the compiler arguments required to compile it. Typically, there is more than one kernel which can perform similar operations. However, each kernel has a unique set of constraints and may result in different performance due to differing code optimizations and input dimensions of the problem. For example, one kernel might be the best choice for large image sizes while another may perform better for smaller image sizes, each using different coding patterns for optimum performance.

All this information is grouped in MIOpen classes collectively called *solvers*. These classes together solve for the best convolution kernel given a problem description. This construct creates a layer of abstraction between the rest of the MIOpen library and kernel specific details, thus all the details of a kernel are completely localized. A solver is trivially constructible by design and therefore has no state, this ensures that kernel compilation launches do not have side effects.

If a developer wishes to add a new kernel to the MIOpen, all that is required is to add the source code for the kernel and implement the associated solver, thereafter the kernel may be selected automatically.

B. Auto tuning infrastructure

In general, any high-performance code leverages auto-tuning for choosing the parameters that may change with the underlying architecture as well as the problem description thereby, impacting performance. MIOpen is not an exception to this rule. This requires that all tunable kernels be tuned for known configuration to achieve maximum performance. Once known, these tuning parameters can be shipped with MIOpen or, the user may employ the same infrastructure to tune MIOpen kernels for custom configurations.

A solver encapsulates the constraints for the tuning parameters as well as the interface machinery to launch tuning instances. The tuning parameters create a grid of possible values of the kernel tuning parameters and the tuning infrastructure compiles and launches a unique kernel for each of these combinations using a pruned search space approach. Once a kernel is tuned and the optimum tuning parameters are known, they are serialized to a designated directory on the user's system for future retrieval.

MIOpen ships with optimized tuning parameters for many configurations used in popular convolutional neural networks. Moreover, the user may run tuning sessions to further optimize kernel codes on their hardware or to add configurations for specific use cases. Further details about the tuning process can be found at [29].

C. Kernel compilation and caching

Launching a kernel requires setting up the compilation parameters and invoking a device-code compiler to generate the binary object. MIOpen device-code consists of kernels written in OpenCL[®], HIP [30] and GCN assembly [31], which may be compiled using *clang* [32].

Since compiling a kernel is a costly and time-consuming procedure, MIOpen employs two levels of caching to improve the runtime performance of the library. This design choice is tightly coupled with how device-code compilers compile and load compute-kernels from the binaries.

Once an kernel file is compiled, it is cached to disk to avoid future compilations of the same source-file with the same parameters. The specific kernel that would be invoked is loaded into memory from the disk and stored in an in-memory cache for subsequent invocation by the same program. This results in substantial runtime improvements since neural network models typically invoke the same kernels many times during an application's lifetime.

Due to the caching effects described above, it is recommended that the user's application performs a *warmup* iteration so that MIOpen's different caches can be populated. Such runs will ensure that subsequent network invocations are accurately timed without the effects of disk I/O or compilation delays. This limitation is not unique to MIOpen and is also applicable to other high-performance libraries.

D. HIP and OpenCL[®] backends

MIOpen supports applications that use the OpenCL[®] and HIP programming models [30]. As shown in Figure 1, all the APIs remain consistent from the client application's perspective, the only difference is in the creation of `miopenHandle` structure, which is created either with a HIP stream or an OpenCL[®] device context. Internally the HIP backend compiles the kernel using an appropriate compiler depending on the kernel source type. Subsequently, the compiled binary object is loaded and passed off to the runtime for execution.

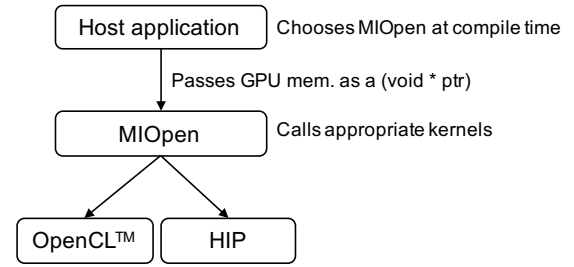


Fig. 1. MIOpen supports OpenCL[®] and HIP Programming Models

IV. MACHINE LEARNING PRIMITIVES

A. Convolution

Most modern neural networks employ convolution as a central operation [33]. Its usefulness and popularity make it a critical piece of the machine learning puzzle, particularly in image processing.

Convolution implementations have a large design space due to its numerical complexity and its diverse inputs make it difficult to generalize across multiple architectures. Over the past few years, different algorithms have been proposed to efficiently compute the convolution of an image with a group of filters. Among these, the Winograd algorithm is notable [34]. The Winograd algorithm achieves the highest efficiency for some key filter sizes. MIOpen's winograd implementation also provides the benefit of not requiring additional "workspace" for intermediate computations. The most general and arguably most expensive in terms of additional storage requirement is to convert the image matrix to a circulant matrix (popularly known as the *im2col* operation [5]), thereafter multiplying the image matrix and the circulant matrix.

Large filter sizes use Fast Fourier Transform (FFT) to convert an image and a filter (after suitable padding) to the frequency domain and then apply a point-wise multiplication followed by the inverse transform to recover the result in time-domain. While this incurs a transformation overhead for the image each time, there are certain cases where this approach is faster than other methods since the filter needs to be transformed only once.

In addition to the above algorithms, MIOpen also implements specialized kernels which directly perform the convolution operation using optimized GCN assembly [31] or OpenCL[®] code [15]. These kernels are collectively known as the *direct* algorithm.

The best performing algorithm is rarely readily apparent on a given architecture for a set of input and filter dimensions. To assess the relative performance of these kernels and return the best performing kernel, MIOpen employs the *find step* before the actual convolution operation. For this step, the user constructs the necessary data structures for the input/output image tensors as well as the convolution descriptor specifying the properties of convolution such as striding, dilation, and padding. The user then calls the MIOpen convolution `Find` API which allows MIOpen to benchmark all the applicable

kernels for the given problem configuration, this information is returned in an array of type `miopenConvAlgoPerf_t`. This enables the library to adjust for any variations in the user hardware and also allows the user to balance the trade-off between execution time and additional memory that may be required for some algorithms.

The `miopenConvAlgoPerf_t` structure mentioned above contains the name of the applicable algorithm, the estimated execution time and the amount of additional memory required by the algorithm. The user may examine this data structure to choose the best algorithm implementation for the problem at hand. This procedure is intended to be performed once and the same find result may be used in subsequent invocations, amortizing its cost.

Types of convolution

Transpose Convolution: Transposed Convolution (also known as deconvolution or fractionally-strided convolution) is an operation typically used to increase the size of the tensor resulting from convolution [35]. The standard convolution operation reduces the size of the image, which is desirable in classification tasks. However, tasks such as image segmentation [36] require the output tensor to have the same size as the input. MIOpen supports transpose convolution required by such networks and may be enabled by setting the `miopenConvolutionMode_t` in `miopenConvolutionDescriptor_t` to `miopenTranspose`.

Depthwise convolution: In depthwise convolution, the input is separated along the depth (channels) and then is convolved with a filter that is also separated along the same axis. The results are stacked into a tensor. To understand why this is useful we have to consider the context in which depthwise convolution occurs – depthwise separable convolution [37]. Depthwise separable convolution involves performing a depthwise convolution followed by a 1x1 convolution on the output tensor called a pointwise convolution [38]. Separating out the process of finding spatial correlation and cross channel correlations, results in fewer parameters as compared to regular convolution [39]. Smaller and more efficient neural networks with depthwise separable convolutions have applications in training on embedded systems such as mobile phones.

Grouped convolutions: Group convolutions were introduced in Alexnet [40], to reduce the memory required for convolution operation. Grouped convolutions are able to achieve accuracy similar to non-grouped convolutions while having fewer parameters. Moreover, grouped convolutions have a higher level of parallelism [41]. Conceptually they are a generalization of depthwise convolution, but instead split the input into individual channels and convolve with a filter that is split up the same way. The results are then stacked together. In grouped convolution the input is split up into groups along the channel axis and is then convolved with a filter that has been grouped along the same axis with the output formed by stacking the resulting tensors from each group; further details may be found in [40].

To perform a groupwise convolution use the function `miopenSetConvolutionGroupCount` on a `miopenConvolutionDescriptor_t` to set the group count. To perform a depthwise convolution use the same function to set group count to the number of channels [38].

Composable Kernels: Different variations of the convolution operation discussed above as well as the variety of algorithms that may be used to implement them make it difficult to develop efficient kernels. One solution to tackle this complexity is to break down these operations into reusable modules that can be universally used by different implementations of different algorithms, and express a kernel as a composition of these modules.

Development work would fall into one of the following categories: 1) describe an algorithm with a hardware-agnostic expression, 2) decide how to map the hardware-agnostic expressions into hardware-dependent modules, 3) implement and optimize the hardware-dependent modules for specific hardware. A potential benefit of breaking down these primitives into smaller modules, is that it then opens new doors to optimization that may fuse these modules together.

This new kernel programming model is referred to as *composable kernels* in MIOpen. MIOpen v2.0 includes an implementation of the implicit GEMM convolution algorithm, using the composable kernel programming approach. Figure 2 gives an overview of the overall structure of direct convolution implementation using this methodology. Further details about this novel programming paradigm will be published in the future.

B. Batch Normalization

Batch normalization is a very successful technique for accelerating deep neural network training. While initially the improved training dynamics were thought to be a result of reduced internal covariate shift, recent research has shown that the true impact of batch normalization layers is a smoother loss function surface making it easier for optimization algorithms to converge to an optimum solution [42], [43].

There are two versions of batch normalization supported in MIOpen: Per-activation and Spatial batch normalization. Per-activation batch normalization is typically positioned after a fully connected layer in a network [42]. For a batch of input samples x represented by channel i , image height and width as j and k respectively, the per activation batch normalization procedure may be described by:

$$y_{ijk} = \gamma_{ijk} \left(\frac{x_{ijk} - \mu_{ijk}}{\sqrt{\sigma_{ijk}^2 + \epsilon}} \right) + \beta_{ijk}$$

Where, y_{ijk} is the output, μ_{ijk} and σ_{ijk}^2 are the mean and variance respectively, ϵ is a small value to avoid division by zero and β_{ijk} and γ_{ijk} are learned parameters.

Batch normalization for convolution layers is termed spatial in that it learns separate parameters γ_i and β_i for each channel, such that the same transform is applied to all

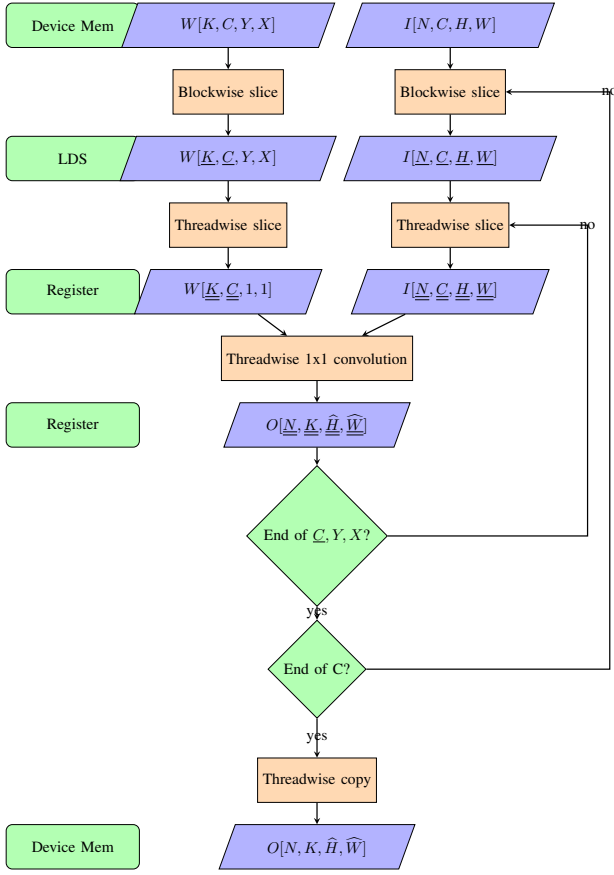


Fig. 2. Work division for convolution using composable kernels

the activations in a single feature map [42]. Likewise, the parameters γ_i and β_i are learned per channel.

Mathematically,

$$y_{ijk} = \gamma_i \left(\frac{x_{ijk} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}} \right) + \beta_i$$

MIOpen supports the batch normalization operation for both training and inference. They all accept the mode parameter from the `miopenBatchNormMode_t` enum, which has two modes `miopenBNPerActivation`, which does element-wise normalization for a fully connected layer and `miopenBNSpatial` which does normalization for convolutions layers. There are specific kernels for training, inference and backward pass for both per activation and spatial batch norm. For more information see [29] and [44].

C. Recurrent Neural Networks

The concept of recurrent neural networks (RNN) dates back to the 1980s for storage of the neuron states in self organizing neural networks [45] [46]. In the Naive RNN structure (also known as vanilla RNN), each neuron is fed with information from input and a previously stored state to predict the next state of the neuron. This attribute of RNNs along with their flexibility in layer construction allows them to substitute for Hidden Markov Model (HMM) to solve state transition

problems such as speech recognition and handwriting recognition [47] [48]. However, vanilla RNN are notoriously difficult to train due to gradient vanishing and exploding issues [9]. As a remedy, long short-term memory (LSTM) structure was introduced which later proved to be effective in sequence learning [47] [48] [49]. Modifications to LSTM such as peephole LSTM, Gated Recurrent Units (GRU) and Simple Recurrent Units (SRU) were later introduced and have been widely used for a variety of applications [50] [51] [52].

MIOpen supports three RNN types prevalent in the industry and research: vanilla RNN, LSTM and, GRU and two kinds of activation function for the hidden state of vanilla RNN neuron: Rectified Linear Unit (ReLU) and hyperbolic tangent (Tanh). Furthermore, information through the RNN may flow in the forward direction (unidirectional RNNs) or both in the forward and backward directions (bidirectional RNNs). MIOpen supports all three RNN types in the unidirectional `miopenRNNunidirection` as well as the bidirectional model `miopenRNNbidirection`. Some RNN layers take input sequences directly from the output of a previous layer while others require a transform to align the intermediate vector dimension or simply to achieve better results. MIOpen satisfies this requirement by supporting two input types: 1) `miopenRNNlinear`, which performs a linear transform before feeding the input to the neuron, and 2) `miopenRNNskip`, which allows a direct input into the neuron. Similarly, bias to the neural network may be added or removed by choosing the mode `miopenRNNWithBias` or `miopenRNNNoBias`.

The dependence of current state on the previous state as well as different RNN configurations make it difficult to achieve high computational efficiency on a GPU platform. Prevalent frameworks such as TensorFlow encapsulate the state updating functions of the RNN neuron in a cell format to achieve better compatibility in different modes, though the impact of the data layouts and computation procedures on performance is neglected [53]. MIOpen handles the RNN computation by taking advantage of two powerful ROCm platform GEMM libraries (1) rocBLAS for the HIP backend, and (2) MIOpenGEMM for the OpenCL® backend, which are augmented by specialized MIOpen kernels for other primitive functions.

1) *Fusion and LSTM*: In the following paragraphs, details about the fusion and optimization of LSTM's forward and backward paths are presented. Similar ideas extend to the design of vanilla RNN as well as GRU in MIOpen. The interested reader is referred to the MIOpen code repository for further details [44].

LSTM has four gates to support its storage and update in long and short-term memories. The structure of LSTM neuron is shown in Figure 3 and its recurrent logic are shown in equation 1-10.

$$s_i t = W_i x_t + R_i h_{t-1} \quad (1)$$

$$s_f t = W_f x_t + R_f h_{t-1} \quad (2)$$

$$s_o t = W_o x_t + R_o h_{t-1} \quad (3)$$

$$s\tilde{c}_t = W_c x_t + R_c h_{t-1} \quad (4)$$

In the equations, the subscript t denotes the time index in the sequence, and i_t , f_t , o_t and c_t denote the updates at the *input*, *forget*, *output* and *cell* gates respectively at time t . The input is denoted as x , and h indicates the hidden state of the LSTM cell. The matrices W_i , W_f , W_o , W_c are the weight matrices of x for the input, forget, output and cell gates respectively. Similarly, the R matrices with appropriate subscripts are the gain matrices for the hidden state h . The above equations represent the linear transformation of the various LSTM states, these interim states go through different activation functions as follows:

$$i_t = \text{sigmoid}(s i_t) \quad (5)$$

$$f_t = \text{sigmoid}(s f_t) \quad (6)$$

$$o_t = \text{sigmoid}(s o_t) \quad (7)$$

$$\tilde{c}_t = \tanh(s\tilde{c}_t) \quad (8)$$

Finally, the cell state and the hidden state for the current time step are calculated as:

$$c_t = f_t \times c_{t-1} + i_t \times \tilde{c}_t \quad (9)$$

$$h_t = o_t \times \tanh(c_t) \quad (10)$$

Note that, in the above equations, each state will be updated by the same input vector x_t and hidden state vector h_{t-1} . Equation 1-4 can then be represented by a single GEMM as described in equation 11:

$$\begin{bmatrix} s i_t \\ s f_t \\ s o_t \\ s\tilde{c}_t \end{bmatrix} = \begin{bmatrix} W_i \\ W_f \\ W_o \\ W_c \end{bmatrix} x_t + \begin{bmatrix} R_i \\ R_f \\ R_o \\ R_c \end{bmatrix} h_{t-1} \quad (11)$$

where $s i_t$, $s f_t$, $s o_t$, $s\tilde{c}_t$ form a large tensor constructed by concatenating the individual buffers. Since the input vectors at different time steps are independent of each other, the computations for all time steps can be further fused in a single GEMM call as illustrated in equation 12 below.

$$\begin{bmatrix} s_0 & s_1 & \dots & s_{T-1} \end{bmatrix} = W \begin{bmatrix} x_0 & x_1 & \dots & x_{T-1} \end{bmatrix} \quad (12)$$

where

$$s_t = \begin{bmatrix} s i_t \\ s f_t \\ s o_t \\ s\tilde{c}_t \end{bmatrix} \quad (13)$$

and

$$W = \begin{bmatrix} W_i \\ W_f \\ W_o \\ W_c \end{bmatrix} \quad (14)$$

For an input sequence of length T , the above optimization yields the following advantages (1) input weight matrices W for all four gates only need to be loaded once, for x_t over all

time steps, leading to (T-1) savings in memory reading of the four matrices; (2) at each time step, the number of loads for the hidden state vector at the previous time step h_{t-1} in each neuron is reduced from four to one. Meanwhile, the operations in equations 5-7 are also fused into one call of the sigmoid kernel due to the computational homogeneity and contiguous memory-layout.

The backward path illustrated in Figure 4 adopts a similar optimization strategy. After deriving the back-propagation error $\Delta s i_t$, $\Delta s f_t$, $\Delta s o_t$, $\Delta s\tilde{c}_t$ at all four gates, the error propagating to the previous state can be derived as depicted in equation 15:

$$\Delta h_{t-1} = \begin{bmatrix} R_i^T & R_f^T & R_o^T & R_c^T \end{bmatrix} \begin{bmatrix} \Delta s i_t \\ \Delta s f_t \\ \Delta s o_t \\ \Delta s\tilde{c}_t \end{bmatrix} + \Delta y_{t-1} \quad (15)$$

while Δy_{t-1} is the error propagated from the higher stack of LSTM layer and can be populated in Δh_{t-1} buffer beforehand. After updating state errors of each gate over all time, a single GEMM call yields the back-propagation error Δx_t for the lower LSTM layer as shown in equations 16 and 17.

$$\Delta x_t = \begin{bmatrix} W_i^T & W_f^T & W_o^T & W_c^T \end{bmatrix} \begin{bmatrix} \Delta s i_t \\ \Delta s f_t \\ \Delta s o_t \\ \Delta s\tilde{c}_t \end{bmatrix} \quad (16)$$

$$\begin{bmatrix} \Delta x_0 & \Delta x_1 & \dots & \Delta x_{T-1} \end{bmatrix} = W^T \begin{bmatrix} \Delta s_0 & \Delta s_1 & \dots & \Delta s_{T-1} \end{bmatrix} \quad (17)$$

In the ideal case, the backpropagation error for weights may be updated using two GEMM calls if the batch sizes are the same for all time steps. The input weight update at each time step is given by:

$$\begin{bmatrix} \Delta W_i & \Delta W_f & \Delta W_o & \Delta W_c \end{bmatrix} = \begin{bmatrix} \Delta s i_t & \Delta s f_t & \Delta s o_t & \Delta s\tilde{c}_t \end{bmatrix} x_t^T \quad (18)$$

In MIOpen, the input weight update over all time steps is given by a single GEMM call as in equation 19:

$$\Delta W = \begin{bmatrix} \Delta s_0 & \Delta s_1 & \dots & \Delta s_{T-1} \end{bmatrix} \begin{bmatrix} x_0^T \\ x_1^T \\ \dots \\ x_{T-1}^T \end{bmatrix} \quad (19)$$

The hidden state weight update at each time step is given by:

$$\begin{bmatrix} \Delta R_i & \Delta R_f & \Delta R_o & \Delta R_c \end{bmatrix} = \begin{bmatrix} \Delta s i_t & \Delta s f_t & \Delta s o_t & \Delta s\tilde{c}_t \end{bmatrix} h_{t-1}^T \quad (20)$$

Similar to the single-GEMM call for input weight update, the hidden state weight update over all time is

$$\Delta R = \begin{bmatrix} \Delta s_0 & \Delta s_1 & \dots & \Delta s_{T-1} \end{bmatrix} \begin{bmatrix} h_{-1}^T \\ h_0^T \\ \dots \\ h_{T-2}^T \end{bmatrix} \quad (21)$$

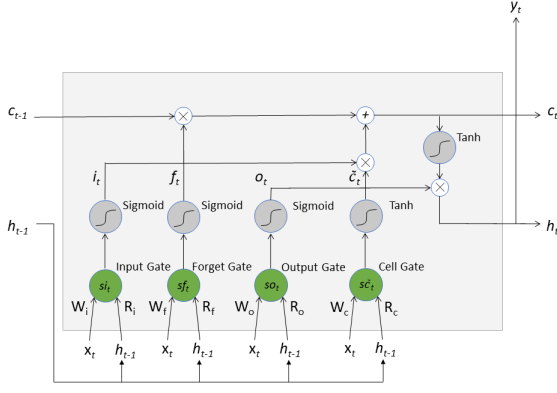


Fig. 3. structure of LSTM neuron and forward flow

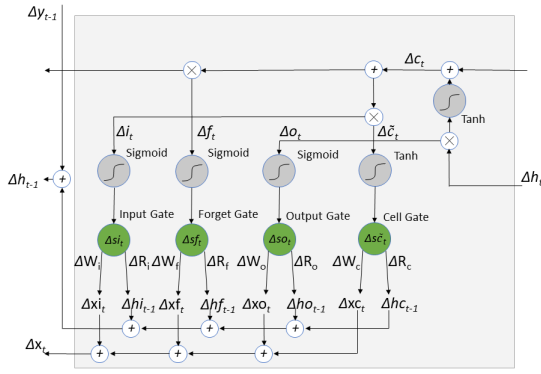


Fig. 4. backward flow of LSTM

In parallel computing, data tensors are usually packed in batches. However, when training different lengths of sentences in an LSTM model, the batch size at each time step can be different. MIOpen requires a length-descending arrangement for batched sentences (longest sentence at the top of the batch while the shortest at the bottom) to guarantee computational efficiency. In practice, consistent batch size along time axis is preferred to achieve the best performance. For instance, in backward weight update, if the batch size varies along the time axis, the GEMM results of hidden state vectors will have to be aligned at each time step and then subsequently accumulated. This will result in $T + 1$ separate GEMM calls in total, instead of just 2.

D. Other Primitives

In addition to the operations above, other operations are required to support the bulk of the computational workload in popular neural network architectures. Among these operations MIOpen implements the following operations for both training and inference:

- 1) Activation Operations
- 2) Pooling

- 3) Softmax
- 4) CTC Loss Function
- 5) Tensor Operators
- 6) Local response normalization

The procedure to invoke these operations is similar to convolution with the exception that they do not require the *find step*.

V. FUSION API

Most neural networks are data-flow graphs where data flows from one direction and is operated upon as it moves from one layer to another. While conceptually data is flowing only in one direction, the underlying kernels implementing these operations have to read data from the global memory, operate on the data and then write the result back for layers down the pipeline. This is necessary due to the limited on-chip memory of the GPUs given the large image and filter sizes in neural network architectures.

However, not all operations require that data be read from and written back to the global memory each time. That is some operations may be fused to increase the compute efficiency of these kernels. This merger of the operations to be performed by a single kernel may be termed as *kernel-fusion*.

As a simple example let's consider an addition operation followed by a rectified linear unit (ReLU) operation. In this case, the intermediate result need not be written back to the main memory, and both the operations may be performed while the individual data elements are in the on-chip memory. Another common sequence of operations is convolution followed by a bias (addition) and ReLU operation. It must be kept in mind that fusions for other operators are much more involved such as the fusion of the convolution and batch normalization operation.

The MIOpen library offers the fusion API to facilitate the efficient fusion of such operations; it allows the user to specify a sequence of operations that are desired to be fused. Once the user specifies this sequence, MIOpen decides the applicable kernel and compiles it; all this information is encapsulated in the `miopenFusionPlanDescriptor` data structure [44].

If merging of the required fusion sequence is feasible, the compilation step of the fusion plan will return success; thereafter the user would supply the runtime arguments for the kernels such as parameters for different operations. Following which, the user would execute the fusion plan with data pointers for the input and output data. The advantage of separating the compilation step from the argument structure is that the fusion plan which has been compiled once, need not be compiled again for different input values. Figure 5 shows a pictorial representation of the steps required to create a fusion plan. Further details and example code can be found at [29].

A. Metadata graph

Internally MIOpen relies on a constraint specification graph, which when traversed with the attributes of fusion operations results in the applicable kernels. Such a mechanism allows the

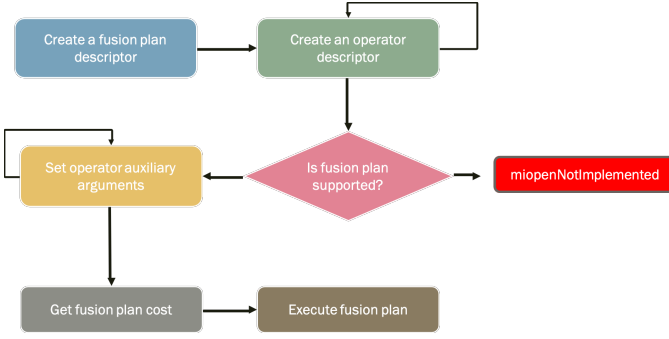


Fig. 5. Steps for creating and executing a fusion plan

addition of new fused kernels with an arbitrary sequence of operations without the combinatorial increase in complexity.

B. Supported Fusions

Tables I and II enumerates the different combinations of fusions that are currently supported by MIOpen for single precision and half precision respectively. The first column indicates the combination of operations that may be fused, where C stands for Convolution, B for bias, N for Batch Normalization and A for activation.

VI. RESULTS

This section highlights the performance improvements that MIOpen is able to offers particularly in convolution as well as some supported fusions. To date, the primary beneficiary of Machine Learning progress has been machine vision as well as Natural Language processing. In machine vision, the convolution operation is the primary workhorse due to the low number of parameters required to learn as compared to regular neural networks as well as the favorable mathematical properties. However, the parameters associated with the convolution operations in different deep convolution neural networks have changed considerably. The early CNNs employed larger filter sizes to reduce the height and width of the feature maps and simultaneously increase the number of feature maps. For instance, LeNet [33] employed filters of size 5×5 while, Alexnet [40] contained filters of size 5×5 as well as 11×11 . However, recently [54], [55] networks have almost exclusively relied on smaller filter sizes namely only 1×1 and 3×3 coupled with striding to reduce the size of the feature map.

Figure 6 shows the relative speedup of different convolution configurations as compared to MIOpen’s im2col+GEMM implementation. The configurations shown therein have been selected randomly from different popular networks such as GoogLeNet, Inception v3, and Inception v4 [55] for image classification. The y-axis in Figure 6 shows log of the speedup obtained by MIOpen, while the x-axis shows the labels for different configurations. Each label shows, respectively, the filter height, filter width, input channels, image height, image width, output channels, padding (height) and padding (width) separated by a hyphen (-).

Figures 6a, 6c and 6e depict the performance gains for kernels with filter height and width equal to 1 (1×1 convolutions) in the forward, backwards-data and backwards-weights directions respectively. While mathematically 1×1 convolutions may be described as a pure GEMM operation, still MIOpen may provide substantial performance benefit in certain cases. Similarly, Figures 6b, 6d and 6f show the performance benefit attained for non- 1×1 kernels in the forward, backward-data and backward-weights directions respectively.

As mentioned in Section IV MIOpen employs the Winograd algorithm for applicable convolutions while the 1×1 convolutions are primarily serviced by kernels written in GCN ISA. Due to the efficiency of the Winograd algorithm, MIOpen can speed up many 3×3 convolutions, however, on larger filter sizes it is not as effective due to granularity loss. Wherein MIOpen’s other convolutional kernels step in to provide speedup, however, in some cases, this speedup is not substantial. The MIOpen team is continuously working on new algorithms to improve performance in these areas.

Section V describes the MIOpen Fusion API, which allows the user to fuse many arbitrary combinations of operations to reduce memory traffic and provide performance gains. Figures 7a and 7b depict the speedup achieved using the Fusion API.

Figure 7a indicates the speedup achieved by the fused operations versus the same operations performed individually. The amount of speedup achieved varies with different configurations, with some being accelerated to as high as 2.5 times the separate run-times. It may be noted that higher speedup is achieved for kernels with fewer output features (channels) since a larger bias vector results in the memory system being the bottleneck.

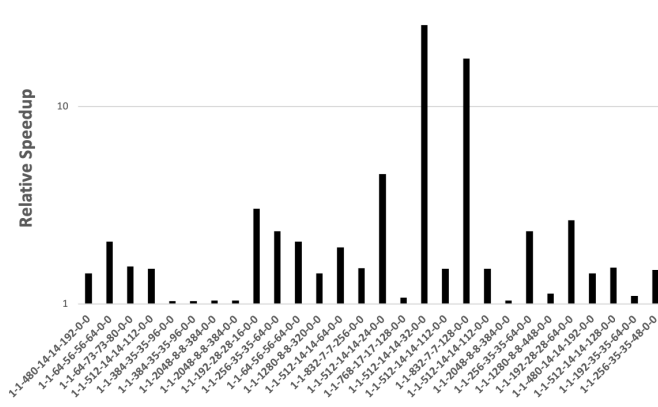
The MIOpen Fusion API is also capable of fusing the Batch-Normalization and Activation operation in the forward direction. The speedup achieved using this fusion for different configurations is depicted in Figure 7b, where the horizontal axis indicates the number of input channels, the height and width of the image. The results indicate that this fusion is more effective for larger image sizes with more number of channels, while smaller images are not able to benefit from the fused operations. However, the possible speedup using this fusion makes it a viable optimization venue to be explored. The MIOpen team is working on expanding the scope of the effectiveness of the Fusion API by enabling more fusions and improving the efficacy of the existing fusions.

VII. CONCLUSIONS AND FUTURE WORK

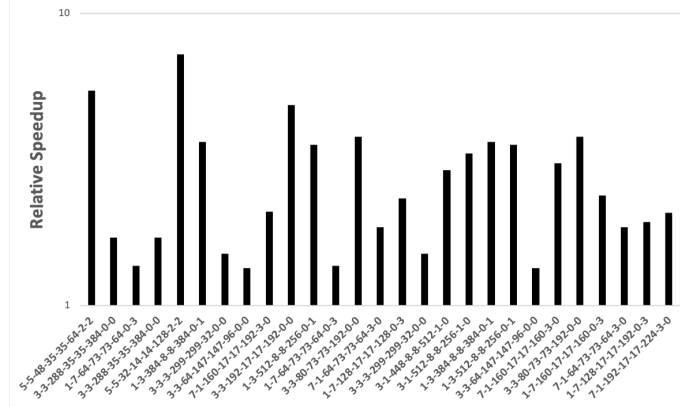
This paper identified some of the challenges faced by a high performance computing library and some of the mechanisms implemented in MIOpen to address these challenges were presented. The open source nature of MIOpen makes it easy for researchers and academics to experiment and implement novel solutions to these problems, the authors look forward to constructive feedback from the community.

ACKNOWLEDGEMENTS

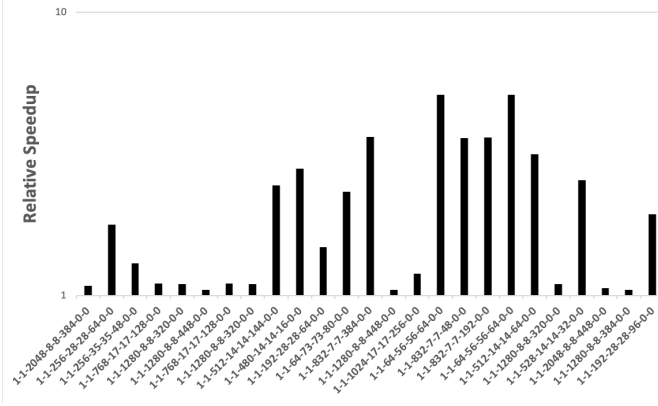
The MIOpen team would like to gratefully acknowledge the valuable contributions of Alex Lyashevsky, James Newling and



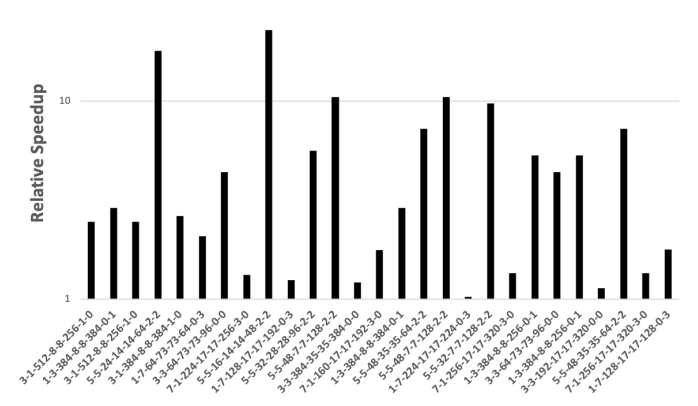
(a) 1x1 filter size (Forward)



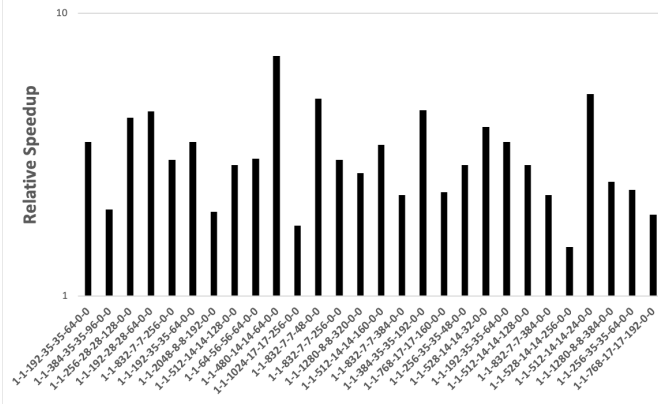
(b) non 1x1 filter sizes (Forward)



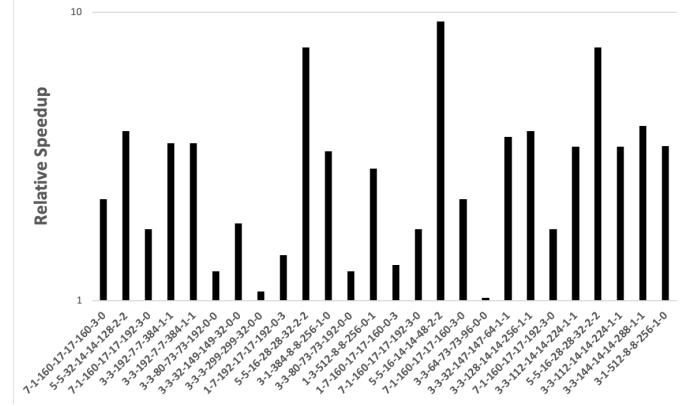
(c) 1x1 filter size (Backward Data)



(d) non 1x1 filter sizes (Backward Data)



(e) 1x1 filter size (Backward Weights)



(f) non 1x1 filter sizes (Backward Weights)

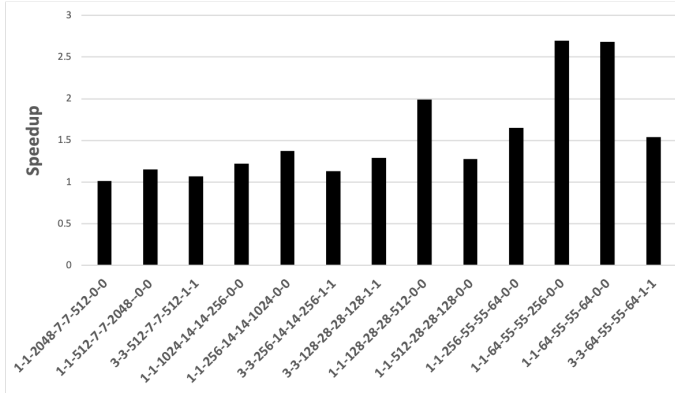
Fig. 6. Relative performance improvement for different convolution configurations as compared to im2col+GEMM

TABLE I
FUSIONS SUPPORTED BY MIOpen (SINGLE PRECISION)

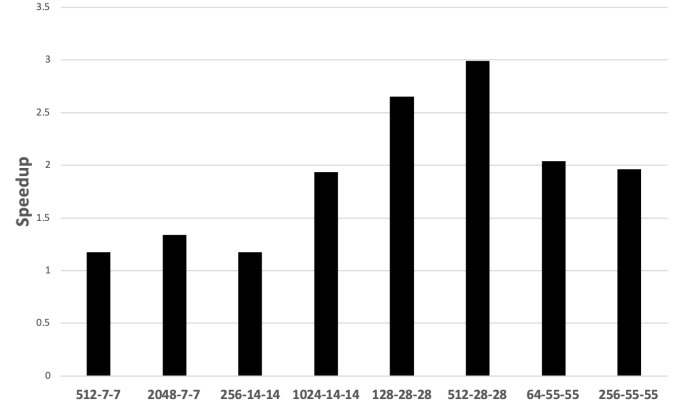
Combination	Conv Algo	Stride	Filter Dims	BN Mode	Activations	Other Constraints
CBNA	Direct	1 and 2	3x3, 5x5, 7x7, 9x9, 11x11	All	All	stride and padding must be either 1 or 2
CBA	Direct		1x1	N/A	All	stride/ padding not supported
		1	1x1, 2x2		Relu, Leaky Relu	c >= 18
		1	3x3		Relu, Leaky Relu	c >= 18 and c is even
		1	4x4, 5x5, 6x6		Relu, Leaky Relu	4 x c >= 18
		1	7x7, 8x8, 9x9		Relu, Leaky Relu	12 x c >= 18
	Winograd	1	10x10, 11x11, 12x12		Relu, Leaky Relu	16 x c >= 18
		1	larger filter sizes		Relu, Leaky Relu	none
		2	1x1		Relu, Leaky Relu	2 x c >= 18
		2	2x2, 3x3, 4x4, 5x5, 6x6		Relu, Leaky Relu	4 x c >= 18
		2	7x7		Relu, Leaky Relu	12 x c >= 18
		2	8x8, 9x9, 10x10, 11x11, 12x12		Relu, Leaky Relu	16 x c >= 18
		2	larger filter sizes		Relu, Leaky Relu	none
NA	-	-	-	All	All	Padding not supported

TABLE II
FUSIONS SUPPORTED BY MIOpen (HALF PRECISION)

Combination	Conv Algo	Stride	Filter Dims	BN Mode	Activations	Other Constraints
CBNA	Direct	1 and 2	3x3, 5x5, 7x7, 9x9, 11x11	All	All	stride and padding must be either 1 or 2
CBA	Direct		1x1		All	stride/ padding not supported



(a) Speedup with Fusing Convolution + Bias + Activation



(b) Speedup with Fusing Batchnorm + Activation

Fig. 7. Relative performance improvement for different fused configurations compared to their non-fused counterparts

the GitHub user `ghostplant` as well as the support of the open source community.

2019 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo and combinations thereof are trademarks of Advanced Micro Devices, Inc. Other product names used in this publication are for identification purposes only and may be trademarks of their respective companies. Python is a trademark of the Python Software Foundation, OpenCL is a trademark of Apple Inc. used by permission by Khronos Group, Inc.

REFERENCES

- [1] Z. Zhang, A. Romero, M. J. Muckley, P. Vincent, L. Yang, and M. Drozdal, "Reducing uncertainty in undersampled MRI reconstruction with active acquisition," *arXiv preprint arXiv:1902.03051*, 2019.
- [2] Y. You, Z. Zhang, C.-J. Hsieh, J. Demmel, and K. Keutzer, "ImageNet training in minutes," in *Proceedings of the 47th International Conference on Parallel Processing*. ACM, 2018, p. 1.
- [3] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard *et al.*, "TensorFlow: A system for large-scale machine learning," in *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, 2016, pp. 265–283.
- [4] A. Paszke, S. Gross, S. Chintala, and G. Chanan, "PyTorch: Tensors and dynamic neural networks in python with strong GPU acceleration," 2017, <https://pytorch.org/>.
- [5] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," *arXiv preprint arXiv:1408.5093*, 2014.
- [6] "MLIR: Multi-level Intermediate Representation for Compiler Infrastructure," <https://github.com/tensorflow/mlir>.
- [7] S. Cyphers, A. K. Bansal, A. Bhiwandiwala, J. Bobba, M. Brookhart, A. Chakraborty, W. Constable, C. Convey, L. Cook, O. Kanawi, R. Kimball, J. Knight, N. Korovaiko, V. Kumar, Y. Lao, C. R. Lishka, J. Menon, J. Myers, S. A. Narayana, A. Procter, and T. J. Webb, "Intel nGraph: An Intermediate Representation, Compiler, and Executor for Deep Learning," pp. 6–8, 2018. [Online]. Available: <http://arxiv.org/abs/1801.08058>
- [8] C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh, "Basic Linear Algebra Subprograms for Fortran usage," 1977.

- [9] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [10] AMD Inc., "ROCm - Open Source Platform for HPC and Ultrascaple GPU Computing," <https://github.com/ROCmSoftwarePlatform>.
- [11] C. Leary and T. Wang, "XLA: TensorFlow, compiled," *TensorFlow Dev Summit*, 2017.
- [12] N. Rotem, J. Fix, S. Abdulrasool, G. Catron, S. Deng, R. Dzhabarov, N. Gibson, J. Hegeman, M. Lele, R. Levenstein *et al.*, "Glow: Graph lowering compiler techniques for neural networks," *arXiv preprint arXiv:1805.00907*, 2018.
- [13] C. Lattner and J. Pienaar, "MLIR Primer: A Compiler Infrastructure for the End of Moores Law," 2019.
- [14] T. Chen, T. Moreau, Z. Jiang, H. Shen, E. Q. Yan, L. Wang, Y. Hu, L. Ceze, C. Guestrin, and A. Krishnamurthy, "TVM: end-to-end optimization stack for deep learning," *arXiv preprint arXiv:1802.04799*, pp. 1–15, 2018.
- [15] A. Munshi, "The OpenCL specification," in *Hot Chips 21 Symposium (HCS)*, 2009 IEEE. IEEE, 2009, pp. 1–314.
- [16] R. C. Whaley and J. J. Dongarra, "Automatically tuned linear algebra software," in *SC'98: Proceedings of the 1998 ACM/IEEE conference on Supercomputing*. IEEE, 1998, pp. 38–38.
- [17] G. Belter, E. R. Jessup, I. Karlin, and J. G. Siek, "Automating the generation of composed linear algebra kernels," in *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*. ACM, 2009, p. 59.
- [18] M. Frigo and S. G. Johnson, "FFTW: An adaptive software architecture for the FFT," in *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP'98 (Cat. No. 98CH36181)*, vol. 3. IEEE, 1998, pp. 1381–1384.
- [19] N. Vasilache, O. Zinenko, T. Theodoridis, P. Goyal, Z. DeVito, W. S. Moses, S. Verdoolaege, A. Adams, and A. Cohen, "Tensor Comprehensions: Framework-Agnostic High-Performance Machine Learning Abstractions," vol. 2, 2018. [Online]. Available: <http://arxiv.org/abs/1802.04730>
- [20] Intel Corporation, "Intel MKL-DNN," <https://software.intel.com/en-us/articles/introducing-dnn-primitives-in-intel-mkl>.
- [21] S. Chetlur, C. Woolley, P. Vandermersch, J. Cohen, J. Tran, B. Catanzaro, and E. Shelhamer, "cuDNN: Efficient Primitives for Deep Learning," pp. 1–9, 2014. [Online]. Available: <http://arxiv.org/abs/1410.0759>
- [22] J. Ragan-Kelley, C. Barnes, A. Adams, S. Paris, F. Durand, and S. Amarasinghe, "Halide: a language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines," *Acm Sigplan Notices*, vol. 48, no. 6, pp. 519–530, 2013.
- [23] F. Kjolstad, S. Kamil, S. Chou, D. Lugato, and S. Amarasinghe, "The tensor algebra compiler," *Proceedings of the ACM on Programming Languages*, vol. 1, no. OOPSLA, p. 77, 2017.
- [24] F. Kjolstad, S. Kamil, J. Ragan-Kelley, D. I. Levin, S. Sueda, D. Chen, E. Vouga, D. M. Kaufman, G. Kanwar, W. Matusik *et al.*, "Simit: A language for physical simulation," *ACM Transactions on Graphics (TOG)*, vol. 35, no. 2, p. 20, 2016.
- [25] M. Luján, T. Freeman, and J. R. Gurd, "OoLaLa: an object oriented analysis and design of numerical linear algebra," in *ACM SIGPLAN Notices*, vol. 35, no. 10. ACM, 2000, pp. 229–252.
- [26] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang, "MXNet: A Flexible and Efficient Machine Learning Library for Heterogeneous Distributed Systems," pp. 1–6, 2015. [Online]. Available: <http://arxiv.org/abs/1512.01274>
- [27] R. Collobert, K. Kavukcuoglu, and C. Farabet, "Implementing neural networks efficiently," in *Neural Networks: Tricks of the Trade*. Springer, 2012, pp. 537–557.
- [28] R. Al-Rfou, G. Alain, A. Almahairi, C. Angermueller, D. Bahdanau, N. Ballas, F. Bastien, F. Bayer, A. Belikov, A. Belopolsky *et al.*, "Theano: A python framework for fast computation of mathematical expressions," *arXiv preprint arXiv:1605.02688*, 2016.
- [29] "MIOpen: Documentation," <https://rocmsoftwareplatform.github.io/MIOpen/doc/html>.
- [30] "AMD HIP," <https://github.com/ROCm-Developer-Tools/HIP>.
- [31] "AMD GCN ISA," <https://developer.amd.com/resources/developer-guides-manuals>.
- [32] "Clang: a C language family frontend for LLVM," <http://clang.llvm.org/>.
- [33] Y. LeCun, Y. Bengio *et al.*, "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, vol. 3361, no. 10, p. 1995, 1995.
- [34] A. Lavin and S. Gray, "Fast Algorithms for Convolutional Neural Networks," 2015. [Online]. Available: <http://arxiv.org/abs/1509.09308>
- [35] V. Dumoulin and F. Visin, "A guide to convolution arithmetic for deep learning," *arXiv preprint arXiv:1603.07285*, 2016.
- [36] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
- [37] L. Sifre and S. Mallat, "Rigid-motion scattering for image classification," *PhD thesis, Ph. D. thesis*, vol. 1, p. 3, 2014.
- [38] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [39] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.
- [40] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [41] Y. Ioannou, D. Robertson, R. Cipolla, and A. Criminisi, "Deep roots: Improving CNN efficiency with hierarchical filter groups," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1231–1240.
- [42] S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.
- [43] H. Shimodaira, "Improving predictive inference under covariate shift by weighting the log-likelihood function," *Journal of statistical planning and inference*, vol. 90, no. 2, pp. 227–244, 2000.
- [44] AMD Inc., "MIOpen: AMD's library for high performance machine learning primitives," <https://github.com/ROCmSoftwarePlatform/MIOpen>.
- [45] J. J. Hopfield, "Neural networks and physical systems with emergent collective computational abilities," *Proceedings of the national academy of sciences*, vol. 79, no. 8, pp. 2554–2558, 1982.
- [46] R. J. Williams, G. E. Hinton, and D. E. Rumelhart, "Learning representations by back-propagating errors," *Nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [47] A. Graves, M. Liwicki, S. Fernandez, R. Bertolami, H. Bunke, and J. Schmidhuber, "A novel connectionist system for unconstrained handwriting recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 31, no. 5, pp. 855–868, 2008.
- [48] A. Graves, A. rahman Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," *IEEE international conference on acoustics, speech and signal processing*, pp. 6645–6649, 2013.
- [49] H. Sak, A. Senior, and F. Beaufays, "Long short-term memory recurrent neural network architectures for large scale acoustic modeling," *Fifteenth annual conference of the international speech communication association*, 2014.
- [50] F. A. Gers and J. Schmidhuber, "Recurrent nets that time and count," *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks. IJCNN 2000. Neural Computing: New Challenges and Perspectives for the New Millennium*, vol. 3, pp. 189–194, 2000.
- [51] K. Cho, B. V. Merriboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.
- [52] T. Lei, Y. Zhang, S. I. Wang, H. Dai, and Y. Artzi, "Simple recurrent units for highly parallelizable recurrence," *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 4470–4481, 2018.
- [53] "TensorFlow RNN cell," https://www.tensorflow.org/api_docs/python/tf/nn/rnn_cell/RNNCell.
- [54] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [55] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.