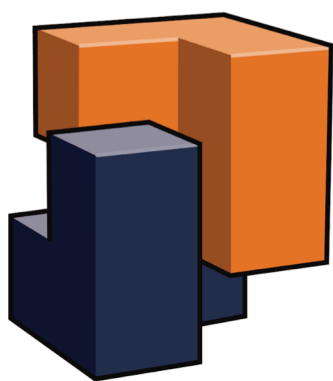


# Bring Your Own Datatypes:

## Enabling Datatype Exploration in Deep Learning with

Gus Smith, Luis Vega, Tianqi Chen, Thierry Moreau, Luis Ceze

CRISP-T3: Scaling Applications and Making the Programmer's Life Easy



# CRISP

Center for Research on Intelligent Storage and Processing in Memory

2780.014 – Programming Framework

### Background

## What is tvm?

An **open-source compiler for deep learning**, easily hackable for research and exploration!



Relay: High-Level Differentiable IR

Tensor Expression IR

LLVM, CUDA, Metal

VTA



Edge  
FPGA

Cloud  
FPGA

ASIC

Optimization

AutoTVM

Hardware  
Fleet

**Relay** is TVM's high-level, differentiable intermediate representation, where we will allow users to insert custom datatypes.

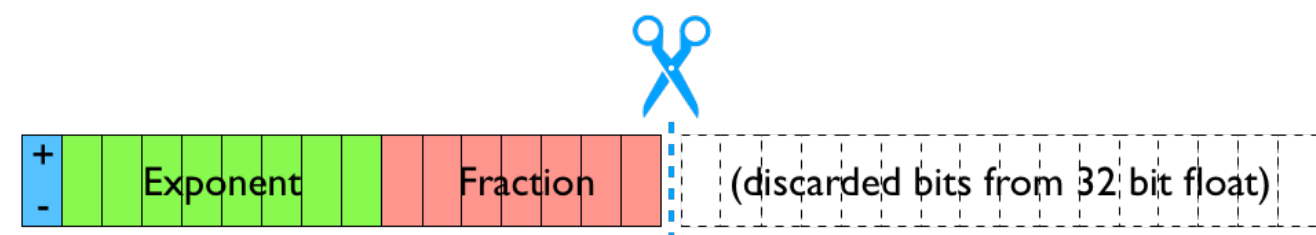
Learn more at [tvm.ai](http://tvm.ai)!

### Background

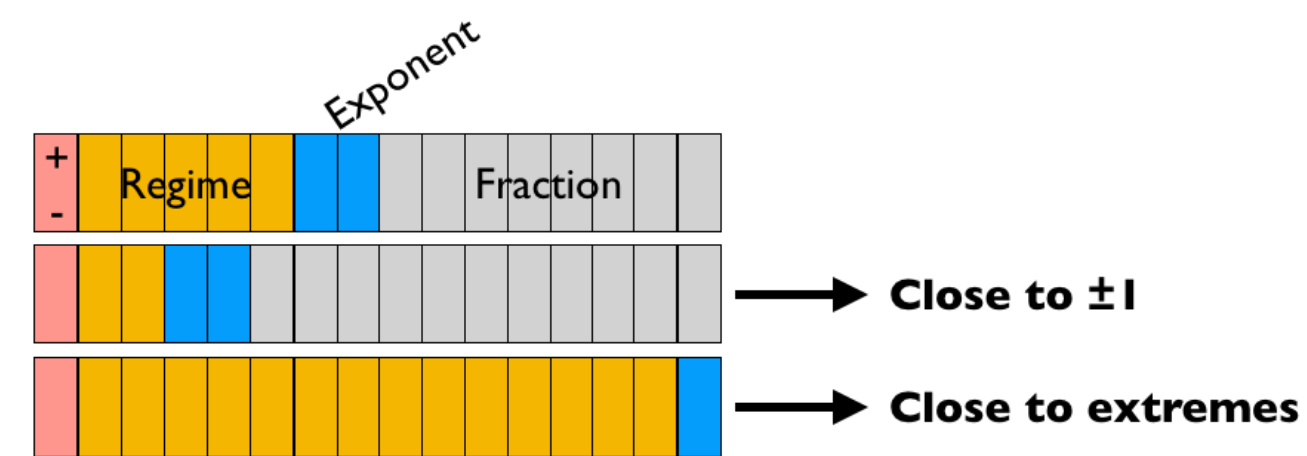
## New Datatypes

To build highly optimized deep learning accelerators, architects are exploring **alternatives to IEEE 754 floating point** for representing numbers in hardware.

One example is the **bfloat16**, which saves space by dropping 16 bits from a float:

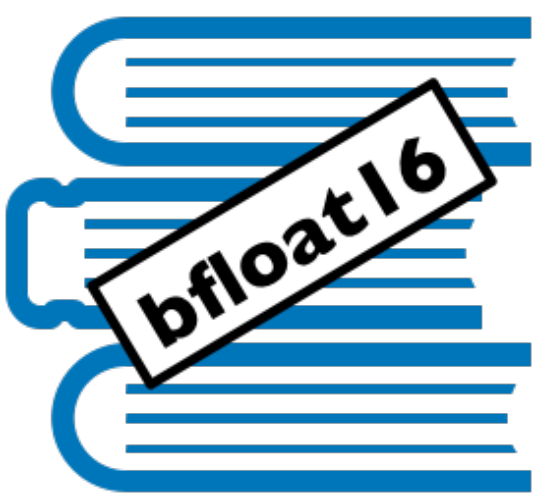


Another example is the **posit**, which uses a variable-size fraction to represent more useful numbers using fewer bits:



### What We Did

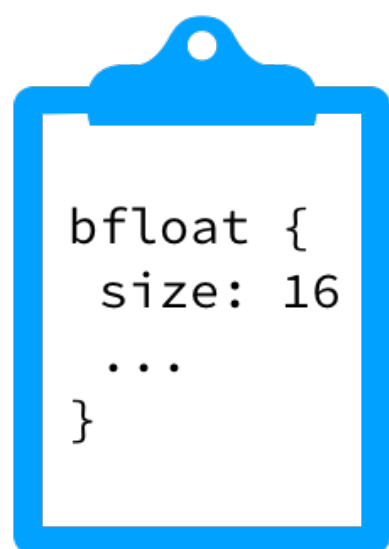
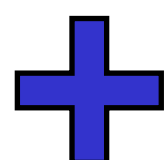
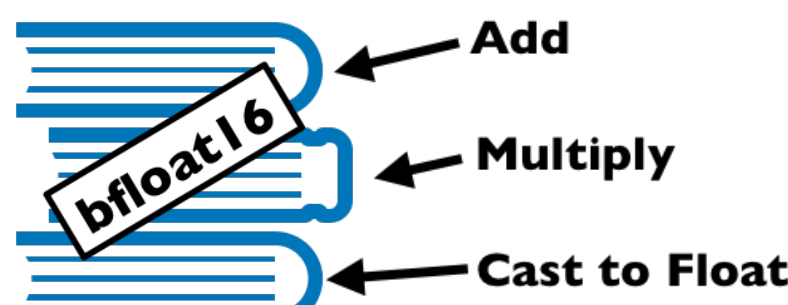
Researchers often prototype new hardware datatypes by emulating them with a **software library**.



We built the **Bring Your Own Datatypes framework** into TVM, which allows users to compile real deep learning workloads which use these emulated datatypes.

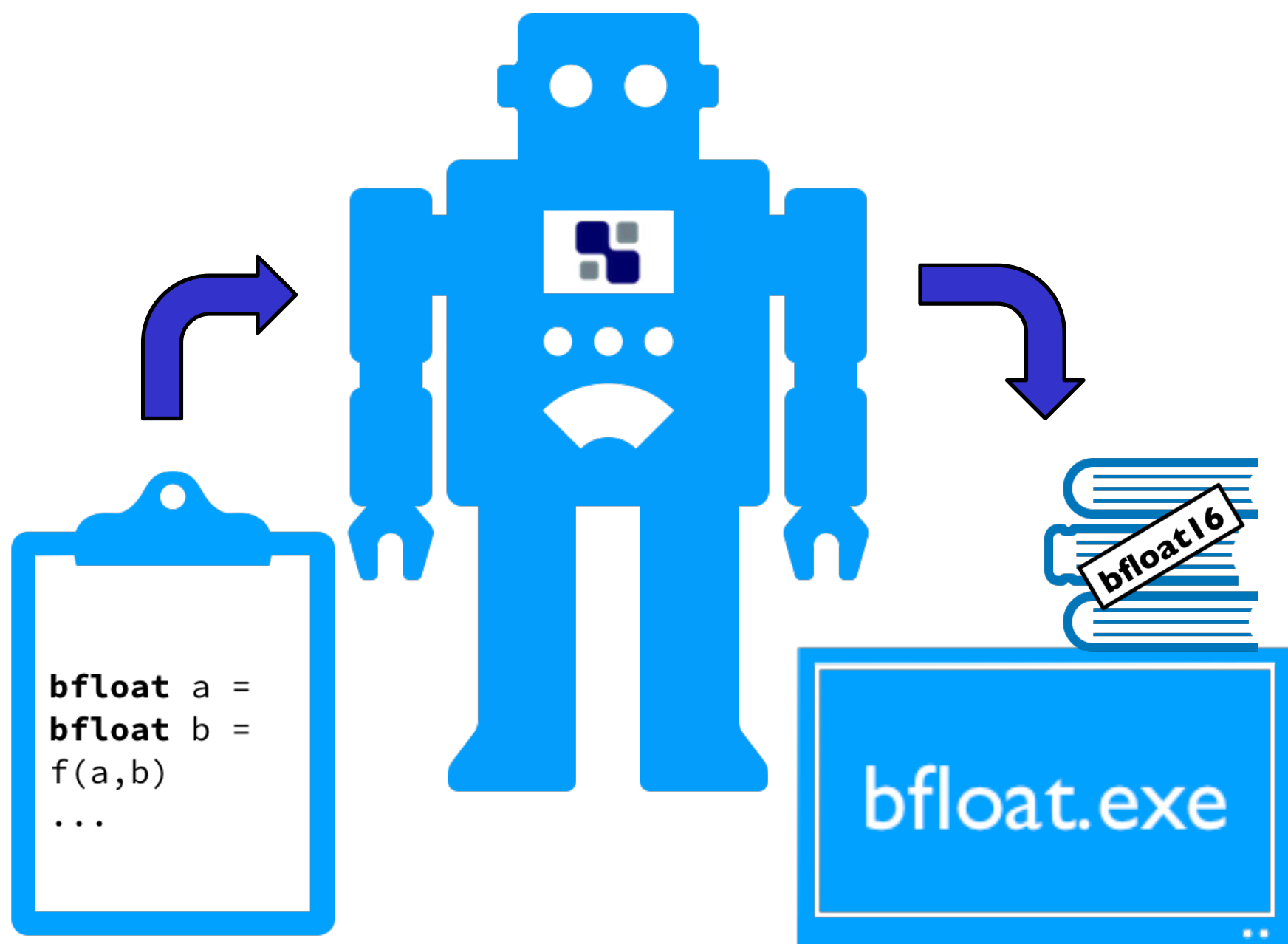
To enable TVM to compile machine learning workloads which utilize custom datatypes, users provide some additional information to the compiler:

First, users **indicate which library functions implement which datatype operators**:



Second, the user **specifies some additional information about the datatype**, such as its size.

With this minimal information, TVM can compile deep learning workloads which rely on custom datatypes!



### Example

```
1 import tvm
2 from tvm import relay
3 from ctypes import RTLD_GLOBAL, CDLL
4
5 # Load the bfloat library into the address space
6 CDLL("libbfloat16.so", RTLD_GLOBAL)
7
8 # Register the datatype with type code 131
9 tvm.datatype.register("bfloat", 131)
10
11 # Register the "Add" operation for bfloats.
12 # The second, third, and fourth arguments specify which operation we are
13 # registering: in this case, an Add of two bfloats when compiling for LLVM.
14 # The first argument tells TVM how to actually lower the bfloat adds:
15 # specifically, it creates a function which turns the bfloat adds into calls to
16 # the BFloat16Add function in libbfloat16.
17 tvm.datatype.register_op(tvm.datatype.create_lower_func("BFloat16Add"), "Add",
18                          "llvm", "bfloat")
19
20 # Create a simple Relay program: Add two bfloats!
21 shape = (3, )
22 dtype = 'custom[bfloat]16'
23 a = relay.var('a', shape=shape, dtype=dtype)
24 b = relay.var('b', shape=shape, dtype=dtype)
25 expr = a + b
26 expr = relay.Function([a, b], expr)
27
28 a_data = ...
29 b_data = ...
30 executor = relay.create_executor()
31
32 # TVM will run the program, calling out to libbfloat16 when it encounters the
33 # bfloat Add!
34 result = executor.evaluate(expr)(a_data, b_data)
```

### Future Work

## What's Next?

Make it **faster**

**Automate** exploration of datatypes in deep learning workloads

Support actual **hardware implementations** of datatypes

### For more info...

For more information, please see the links to the Bring Your Own Datatype slides, presentation, and Python notebook here:

[sampl.cs.washington.edu/tvmfrcr/#program](http://sampl.cs.washington.edu/tvmfrcr/#program)