Bring Your Own Datatypes into TVM

Gus Smith Advisor: Luis Ceze University of Washington







CRISP Center for Research in Intelligent

Storage and Processing in Memory

PAUL G. ALLEN SCHOOL of computer science & engineering



































* * * * * * * * * * * * * * * * * * *



Deep learning needs extreme hardware specialization, and hardware specialization needs new datatypes.







What do we mean by datatypes?



What do we mean by datatypes?

Why do we need new datatypes?



What do we mean by datatypes?

Why do we need new datatypes?

What do new datatypes look like?









An extensible, optimizing compiler for deep learning.

What is TVM?





An extensible, optimizing compiler for deep learning.

What is TVM?

See tvm.ai for more info!









The TVM Stack



High-Level Differentiable IR

Tensor Expression IR





High-Level Differentiable IR

Tensor Expression IR





High-Level Differentiable IR

Tensor Expression IR

LLVM, CUDA, Metal





LLVM, CUDA, Metal













- X = tvm.placeholder(shape=(3,))
- Y = tvm.placeholder(shape=(3,))



X = tvm.placeholder(shape=(3,))

$$Z = topi.add(X, Y)$$



X = tvm.placeholder(shape=(3,))

$$Z = topi.add(X, Y)$$



- X = tvm.placeholder(shape=(3,))
- Y = tvm.placeholder(shape=(3,))
- Z = topi.add(X, Y)

schedule = tvm.create_schedule([Z.op])



X = tvm.placeholder(shape=(3,))

Z = topi.add(X, Y)

schedule = tvm.create_schedule([Z.op]) lowered_function = tvm.lower(schedule, [X, Y, Z])



X = tvm.placeholder(shape=(3,))

Z = topi.add(X, Y)

schedule = tvm.create_schedule([Z.op]) lowered_function = tvm.lower(schedule, [X, Y, Z]) built_program = tvm.build(lowered_function)



X = tvm.placeholder(shape=(3,))

Z = topi.add(X, Y)

schedule = tvm.create_schedule([Z.op]) lowered_function = tvm.lower(schedule, [X, Y, Z]) built_program = tvm.build(lowered_function)



X = tvm.placeholder(shape=(3,))

Z = topi.add(X, Y)

schedule = tvm.create_schedule([Z.op]) lowered_function = tvm.lower(schedule, [X, Y, Z]) built_program = tvm.build(lowered_function)

x = tvm.nd.array(np.random.rand(3).astype("float32")) y = tvm.nd.array(np.random.rand(3).astype("float32"))



X = tvm.placeholder(shape=(3,))

Z = topi.add(X, Y)

schedule = tvm.create_schedule([Z.op]) <u>lowered_function</u> = tvm.lower(schedule, [X, Y, Z]) built_program = tvm.build(lowered_function)

<u>x</u> = tvm.nd.array(np.random.rand(3).astype("float32")) y = tvm.nd.array(np.random.rand(3).astype("float32")) $\underline{z} = tvm.nd.empty((3,))$



X = tvm.placeholder(shape=(3,))

Z = topi.add(X, Y)

schedule = tvm.create_schedule([Z.op]) <u>lowered_function</u> = tvm.lower(schedule, [X, Y, Z]) built_program = tvm.build(lowered_function)

<u>x</u> = tvm.nd.array(np.random.rand(3).astype("float32")) y = tvm.nd.array(np.random.rand(3).astype("float32")) $\underline{z} = tvm.nd.empty((3,))$



X = tvm.placeholder(shape=(3,))

```
Y = tvm.placeholder(shape=(3, ))
```

```
Z = topi.add(X, Y)
```

schedule = tvm.create_schedule([Z.op]) <u>lowered_function</u> = tvm.lower(schedule, [X, Y, Z]) built_program = tvm.build(lowered_function)

<u>x</u> = tvm.nd.array(np.random.rand(3).astype("float32")) y = tvm.nd.array(np.random.rand(3).astype("float32")) $\underline{z} = tvm.nd.empty((3,))$

built_program(x, y, z)



X = tvm.placeholder(shape=(3,))

```
Y = tvm.placeholder(shape=(3, ))
```

```
Z = topi.add(X, Y)
```

schedule = tvm.create_schedule([Z.op]) <u>lowered_function</u> = tvm.lower(schedule, [X, Y, Z]) built_program = tvm.build(lowered_function)

<u>x</u> = tvm.nd.array(np.random.rand(3).astype("float32")) y = tvm.nd.array(np.random.rand(3).astype("float32")) $\underline{z} = tvm.nd.empty((3,))$

built_program(x, y, z)

Vector Add in TVM

X: [0.5727742 0.16838571 0.9647888]

- **y:** [0.75005066 0.12305858 0.9467064]




What do we mean by "datatypes"?

What do we mean by "datatypes"?

What do we mean by "datatypes"?



What do we mean by "datatypes"?





What do we mean by "datatypes"?



What do we mean by "datatypes"?

In the past, this generally meant IEEE 754 floating point:



...and others!

What do we mean by "datatypes"?

Because IEEE floats...

• take up too much space

- take up too much space
- have poor dynamic range

- take up too much space
- have poor dynamic range
- require overly-complex hardware

- take up too much space
- have poor dynamic range
- require overly-complex hardware
- create reproducibility issues

What do new datatypes look like?



What do new datatypes look like?









Greater dynamic range is more useful for deep learning workloads





n = 16, <u>es</u> = 2





another multiplicative factor on the fraction



Posits use a variable-size regime field to add dynamic scaling: regime determines



another multiplicative factor on the fraction

Posits can have large dynamic range even down to 8 bits



Posits use a variable-size regime field to add dynamic scaling: regime determines





Posits use a variable-size **regime** field to add dynamic scaling: **regime** determines another multiplicative factor on the fraction

Posits can have large dynamic range even down to 8 bits

See John Gustafson's "Beating Floating Point at its Own Game"











Block floating point: a group of numbers share an exponent





Block floating point: a group of numbers share an exponent

Plus, they split the fields between host and device!





Block floating point: a group of numbers share an exponent

Plus, they split the fields between host and device!

See Köster et. al., "An Adaptive Numerical Format for Efficient Training of Deep Neural Networks"





Taking inspiration from posits, but with many improvements

Taking inspiration from posits, but with many improvements

E.g. some computation is done in the log space to reduce hardware complexity

- Taking inspiration from posits, but with many improvements
- E.g. some computation is done in the log space to reduce hardware complexity
- Hardware implementation from Jeff Johnson at Facebook is competitive with IEEE floats:

Taking inspiration from posits, but with many improvements

E.g. some computation is done in the log space to reduce hardware complexity

floats:

Component int8/32 MAC PE

- Hardware implementation from Jeff Johnson at Facebook is competitive with IEEE

Area $\mu \mathrm{m}^2$	Power μ W
336.672	283

Taking inspiration from posits, but with many improvements

E.g. some computation is done in the log space to reduce hardware complexity

floats:

Component int8/32 MAC PE (8, 1, 5, 5, 7) log ELMA PE

- Hardware implementation from Jeff Johnson at Facebook is competitive with IEEE

Α	rea $\mu \mathrm{m}^2$	Power μ W
33	36.672	283
37	76.110	272

Taking inspiration from posits, but with many improvements

E.g. some computation is done in the log space to reduce hardware complexity

floats:

Component int8/32 MAC PE (8, 1, 5, 5, 7) log ELMA PE float16 (w/o denormals) FMA PE

- Hardware implementation from Jeff Johnson at Facebook is competitive with IEEE

Α	rea $\mu \mathrm{m}^2$	Power μW
33	36.672	283
37	76.110	272
15	545.012	1358

Taking inspiration from posits, but with many improvements

E.g. some computation is done in the log space to reduce hardware complexity

floats:

Component int8/32 MAC PE (8, 1, 5, 5, 7) log ELMA PE float16 (w/o denormals) FMA PE (5, 10) (11, 11, 10) log ELMA PE

- Hardware implementation from Jeff Johnson at Facebook is competitive with IEEE

Area μm^2	Power μ W
336.672	283
376.110	272
1545.012	1358
1043.154	805

Taking inspiration from posits, but with many improvements

E.g. some computation is done in the log space to reduce hardware complexity

floats:

Component int8/32 MAC PE (8, 1, 5, 5, 7) log ELMA PE float16 (w/o denormals) FMA PE (5, 10) (11, 11, 10) log ELMA PE

See <u>"Rethinking Floating Point for Deep Learning"</u>

- Hardware implementation from Jeff Johnson at Facebook is competitive with IEEE

Area μm^2	Power μ W
336.672	283
376.110	272
1545.012	1358
1043.154	805



Bring Your Own Datatypes


To prototype a hardware datatype, researchers will emulate the datatype in software by building a software library



To prototype a hardware datatype, researchers will emulate the datatype in software by building a software library



To prototype a hardware datatype, researchers will emulate the datatype in software by building a software library



To prototype a hardware datatype, researchers will emulate the datatype in software by building a software library



To prototype a hardware datatype, researchers will emulate the datatype in software by building a software library



To prototype a hardware datatype, researchers will emulate the datatype in software by building a software library



To prototype a hardware datatype, researchers will emulate the datatype in software by building a software library





To prototype a hardware datatype, researchers will emulate the datatype in software by building a software library

To test their datatype, they will hack apart a benchmark or application and shove their software library in

Can we do better? Can we use TVM to compile workloads with new datatypes?









Currently, TVM does not know how to interpret programs with arbitrary datatypes.





Currently, TVM does not know how to interpret programs with arbitrary datatypes.





Currently, TVM does not know how to interpret programs with arbitrary datatypes.

Luckily, TVM is extensible!

























Supporting Datatype Research With TVM Add **Multiply** Cast to Float







Supporting Datatype Research With TVM e.exe my ••







I. User makes or finds a datatype library

- I. User makes or finds a datatype library
- 2. User writes a program using their datatypes directly

- I. User makes or finds a datatype library
- 2. User writes a program using their datatypes directly
- 3. User points TVM to the important functions (+, *) in the library

- I. User makes or finds a datatype library
- 2. User writes a program using their datatypes directly
- 3. User points TVM to the important functions (+, *) in the library
- 4. User gives TVM other information e.g. datatype size

- I. User makes or finds a datatype library
- 2. User writes a program using their datatypes directly
- 3. User points TVM to the important functions (+, *) in the library
- 4. User gives TVM other information e.g. datatype size
- 5. TVM compiles programs, handling the custom datatype by compiling calls into the provided library

Limitations of this Approach

Limitations of this Approach

Currently, we're only supporting software implementations of datatypes.

Limitations of this Approach

the moment, but is planned.

Currently, we're only supporting software implementations of datatypes.

Compiling for hardware implementations of the datatype is out of scope for

Implementation



Datatype Registry



User registers their datatype:

Datatype Registry


User registers their datatype:

tvm.datatype.register("bfloat", 129)

Datatype Registry



User registers their datatype:

tvm.datatype.register("<u>bfloat</u>", 129)

Where 129 is a manually chosen code for the type

Datatype Registry





Then, in the code, we can give tensors the custom datatype:

Then, in the code, we can give tensors the custom datatype: dtype="custom[bfloat]16"



Then, in the code, we can give tensors the custom datatype: dtype="custom[bfloat]16"

Here, "16" is how we tell TVM the size of the datatype.



Operation Registry



Operation Registry



tvm.datatype.register_op(

Operation Registry



tvm.datatype.register_op(lower_func,

Operation Registry



tvm.datatype.register_op(lower_func, "Add", "llvm", "bfloat")

Operation Registry



tvm.datatype.register_op(lower_func, "Add", "llvm", "bfloat")

In the common case, we will use:

Operation Registry



```
tvm.datatype.register_op(
lower_func,
"Add", "llvm", "bfloat")
```

In the common case, we will use:

tvm.datatype.register_op(

Operation Registry



tvm.datatype.register_op(lower_func, "Add", "llvm", "bfloat")

In the common case, we will use:

tvm.datatype.register_op(tvm.datatype.create_lower_func("BFloat16Add"),

Operation Registry



```
tvm.datatype.register_op(
lower_func,
"Add", "llvm", "bfloat")
```

In the common case, we will use:

```
tvm.datatype.register_op(
tvm.datatype.create_lower_func("BFloat16Add"),
"Add", "llvm", "bfloat")
```

Operation Registry



```
tvm.datatype.register_op(
lower_func,
"Add", "llvm", "bfloat")
```

In the common case, we will use:

```
tvm.datatype.register_op(
tvm.datatype.create_lower_func("BFloat16Add"),
"Add", "llvm", "bfloat")
```

Operation Registry



```
tvm.datatype.register_op(
lower_func,
"Add", "llvm", "bfloat")
```

In the common case, we will use:

```
tvm.datatype.register_op(
tvm.datatype.create_lower_func("BFloat16Add"),
"Add", "llvm", "bfloat")
```



Operation Registry



```
tvm.datatype.register_op(
lower_func,
"Add", "llvm", "bfloat")
```

In the common case, we will use:

```
tvm.datatype.register_op(
tvm.datatype.create_lower_func("BFloat16Add"),
"Add", "llvm", "bfloat")
```



Operation Registry



```
tvm.datatype.register_op(
lower_func,
"Add", "llvm", "bfloat")
```

In the common case, we will use:

```
tvm.datatype.register_op(
tvm.datatype.create_lower_func("BFloat16Add"),
"Add", "llvm", "bfloat")
```



Operation Registry









• Short term: evaluating real deep learning models with modern datatypes



- Short term: evaluating real deep learning models with modern datatypes
- Long term: supporting custom datatype hardware









Thank You!

CRISP **Center for Research in Intelligent**

Storage and Processing in Memory



