Bring Your Own Datatypes: **Enabling Number Format Exploration** with TVM Gus Smith University of Washington











• Why you should care about new datatypes



- Why you should care about new datatypes
- How TVM supports datatype exploration



- Why you should care about new datatypes
- How TVM supports datatype exploration Currently: software support only



- Why you should care about new datatypes
- How TVM supports datatype exploration Currently: software support only Future directions: hardware support



- Why you should care about new datatypes
- How TVM supports datatype exploration Currently: **software** support only Future directions: hardware support
- ...and of course, a coding session!





































What's wrong with IEEE floats?



IEEE floats...



IEEE floats...

• take up too much space,



IEEE floats...

- take up too much space,
- have poor dynamic range for their size, and



IEEE floats...

- take up too much space,
- have poor dynamic range for their size, and
- require overly-complex hardware

























Deep learning needs hardware specialization

Deep learning needs hardware specializationand hardware specialization needs new datatypes!

This helps

This helps

• datatype researchers, who need easy ways to test their datatypes,

This helps

- datatype researchers, who need easy ways to test their datatypes,

• but also, all TVM users, who can easily use new datatypes in their models.

Bring Your Own Datatypes


To prototype a hardware datatype, researchers will <u>emulate the datatype</u> <u>in software</u> by building a software library.



To prototype a hardware datatype, researchers will <u>emulate the datatype</u> <u>in software</u> by building a software library.



To prototype a hardware datatype, researchers will <u>emulate the datatype</u> <u>in software</u> by building a software library.



To prototype a hardware datatype, researchers will <u>emulate the datatype</u> <u>in software</u> by building a software library.



To prototype a hardware datatype, researchers will <u>emulate the datatype</u> <u>in software</u> by building a software library.



To prototype a hardware datatype, researchers will <u>emulate the datatype</u> <u>in software</u> by building a software library.





Second State St				
<> Code	Issues 0	្រា Pull requests 0	Projects 0	

		• Watch	3	★ Star	5	% Fork	75,570
🗉 Wiki	C Security	Insights					



Y xman / forked from	tensorflow n tensorflow/tensorfl	ow	
<> Code	Issues 0	1 Pull requests 0	Projects 0
⑦ 36,	645 commits	ှို 14 branch	es
Branch: po	sit 🕶 New pull	request	
This branc	h is 237 commit	s ahead, 22081 commi	ts behind tensorflov
🇊 xman p	oosit: update versio	on	

		• Watch	3	★ Star	5	% Fork	75,570
🗐 Wiki	C Security	Insights					
⇒8 releas	es	41,481 contrib	outors	5	م أ ك	Apache-2	.0
				Find File		lone or do	wnload 🗸
v:master.				រ៉ា Pi	ull req	uest 🖹 (Compare
			Late	est commit	47fa4	c7 on Sep	20, 2018



Y xman / forked from	tensorflow n tensorflow/tensorflo	ow	
<> Code	Issues 0	រឿ Pull requests 0	Projects 0
⑦ 36,	,645 commits	₽ 14 branch	es
Branch: po	sit - New pull r	request	
This brand	ch is 237 commits	s ahead, 22081 commi	ts behind tensorflov
🎆 xman	posit: update versio	n	

Showing 228 changed files with 5,897 additions and 950 deletions.



Unified Split



<pre>% xman / forked from</pre>	tensorflow	low		
<> Code	Issues 0	ឿ Pull requests 0	Projects 0	
@ 36,	645 commits	⊮ 14 branch	es	<
Branch: po	sit 🔻 New pull	request		
This branc	ch is 237 commit	s ahead, 22081 commi	ts behind tensor	flow
🇊 xman p	oosit: update versio	on		

Showing 228 changed files with 5,897 additions and 950 deletions.



Unified Split



Can we do better? Can we let users bring their own datatypes to TVM, and have TVM do the rest of the work?





Out of the box, TVM does not know how to interpret programs with custom datatypes.





Out of the box, TVM does not know how to interpret programs with custom datatypes.





Out of the box, TVM does not know how to interpret programs with custom datatypes.

But, with a bit of information, TVM can compile and run these programs easily!













bfloat { size: 16 • • • Add Multiply Cast to Float



bfloat { size: 16 • • • Add - Multiply Cast to Float



bfloat { size: 16 • • • Add **Multiply** Cast to Float



















I. User makes or finds a library which implements their datatype in software

- I. User makes or finds a library which implements their datatype in software
- 2. User writes a program which uses their custom datatype

- I. User makes or finds a library which implements their datatype in software
- 2. User writes a program which uses their custom datatype
- 3. User points TVM to the important functions (+, *, cast-to-float) in the library

- I. User makes or finds a library which implements their datatype in software
- 2. User writes a program which uses their custom datatype
- 3. User points TVM to the important functions (+, *, cast-to-float) in the library
- 4. User gives TVM other information e.g. datatype size

- I. User makes or finds a library which implements their datatype in software
- 2. User writes a program which uses their custom datatype
- 3. User points TVM to the important functions (+, *, cast-to-float) in the library
- 4. User gives TVM other information e.g. datatype size
- 5. TVM compiles the program, handling the custom datatype by compiling to 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

Compiling for hardware implementations is a work in progress!

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


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



- modern datatypes



• Short term: evaluating real deep learning models with software emulations of

Have a model or datatype you're interested in specifically? I can help!

- Short term: evaluating real deep learning models with software emulations of modern datatypes
 - Have a model or datatype you're interested in specifically? I can help!
- Also short term: improving performance in software (datatype emulation) libraries are slow, and models run many operations!)



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

- Also short term: improving performance in software (datatype emulation) libraries are slow, and models run many operations!)
- **Long term:** supporting custom datatype <u>hardware</u> implementations in VTA



Have a model or datatype you're interested in specifically? I can help!

Let's move to the notebook! (link)







