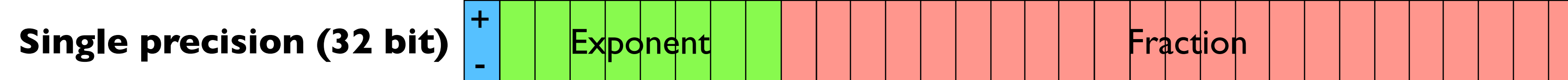


Bring Your Own Datatypes: Enabling Custom Datatype Exploration in Deep Learning

Gus Smith, Qualifying Exam
February 24th, 2020

By “datatypes”, I mean **numerical datatypes**: how the hardware represents and operates on real numbers.

IEEE 754 Floating Point

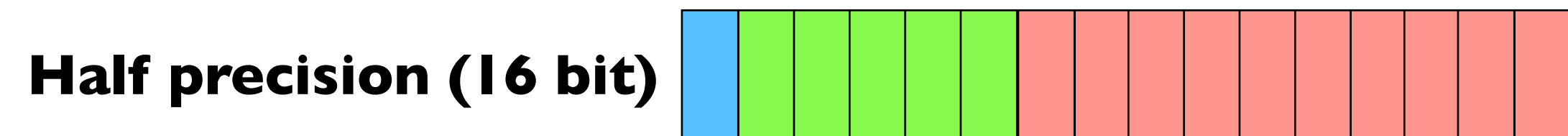


$$\text{value} \approx \text{sign} * 2^{\text{exponent}} * 1.\text{fraction}$$

IEEE 754 Floating Point

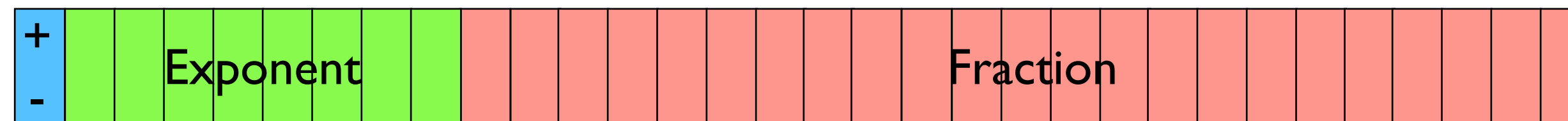


$$\text{value} \approx \text{sign} * 2^{\text{exponent}} * 1.\text{fraction}$$



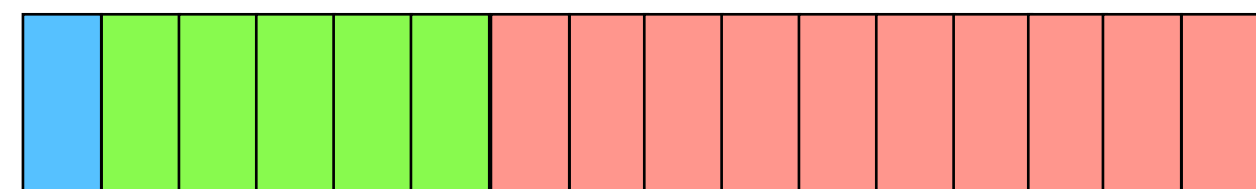
IEEE 754 Floating Point

Single precision (32 bit)

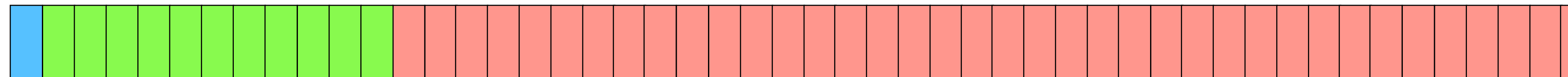


$$\text{value} \approx \text{sign} * 2^{\text{exponent}} * 1.\text{fraction}$$

Half precision (16 bit)



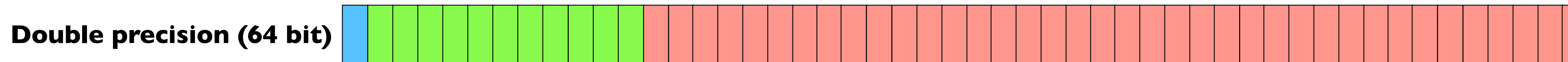
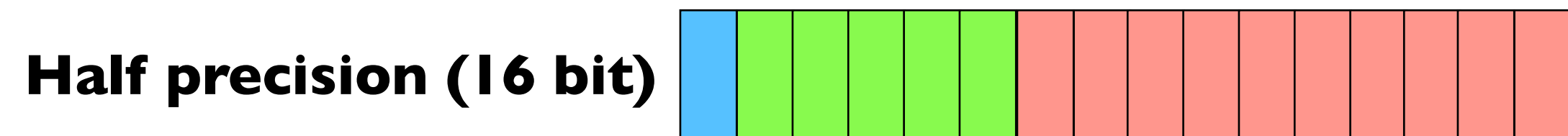
Double precision (64 bit)



IEEE 754 Floating Point

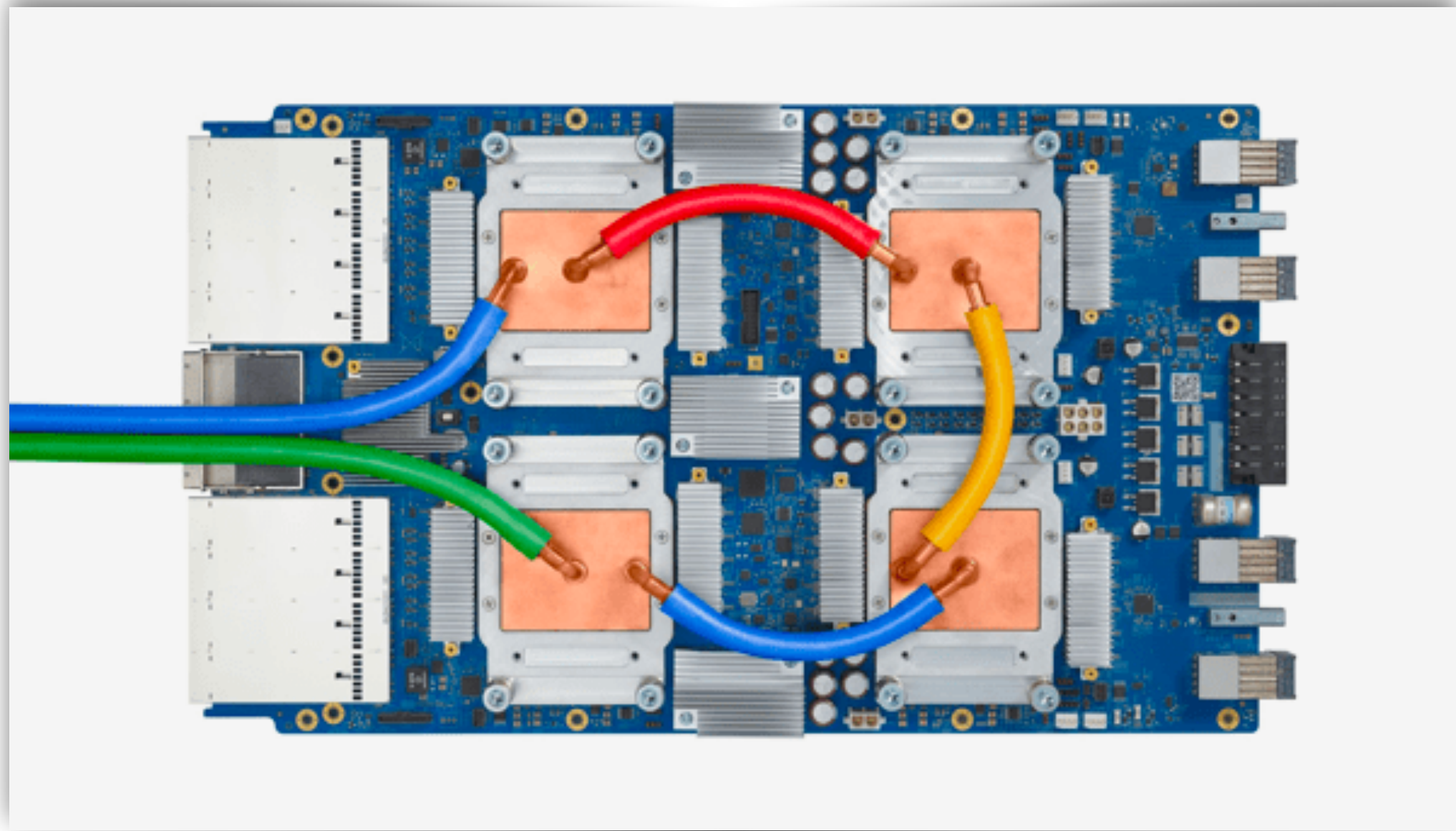


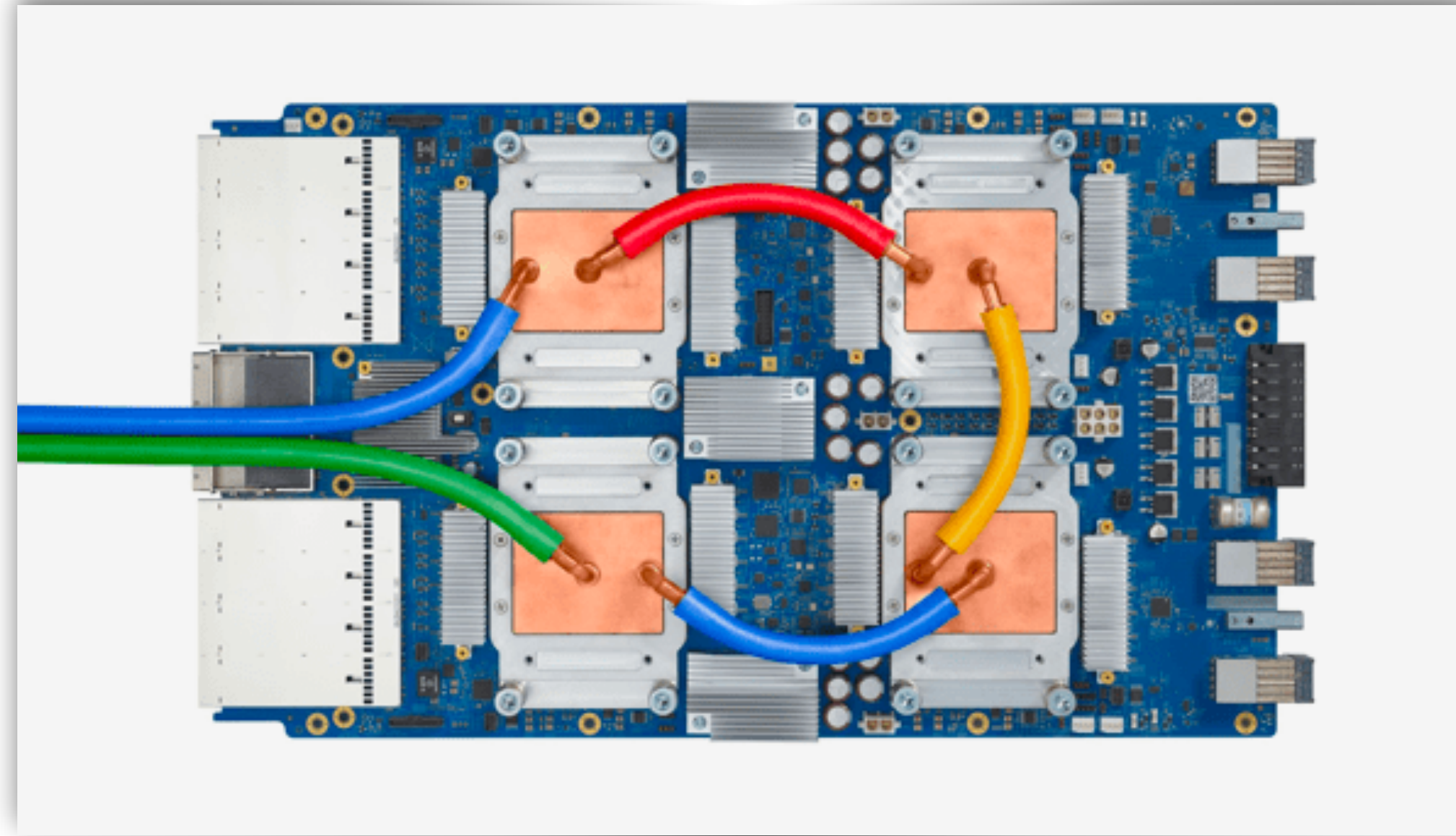
$$\text{value} \approx \text{sign} * 2^{\text{exponent}} * 1.\text{fraction}$$



Has remained an industry standard for more than thirty years!

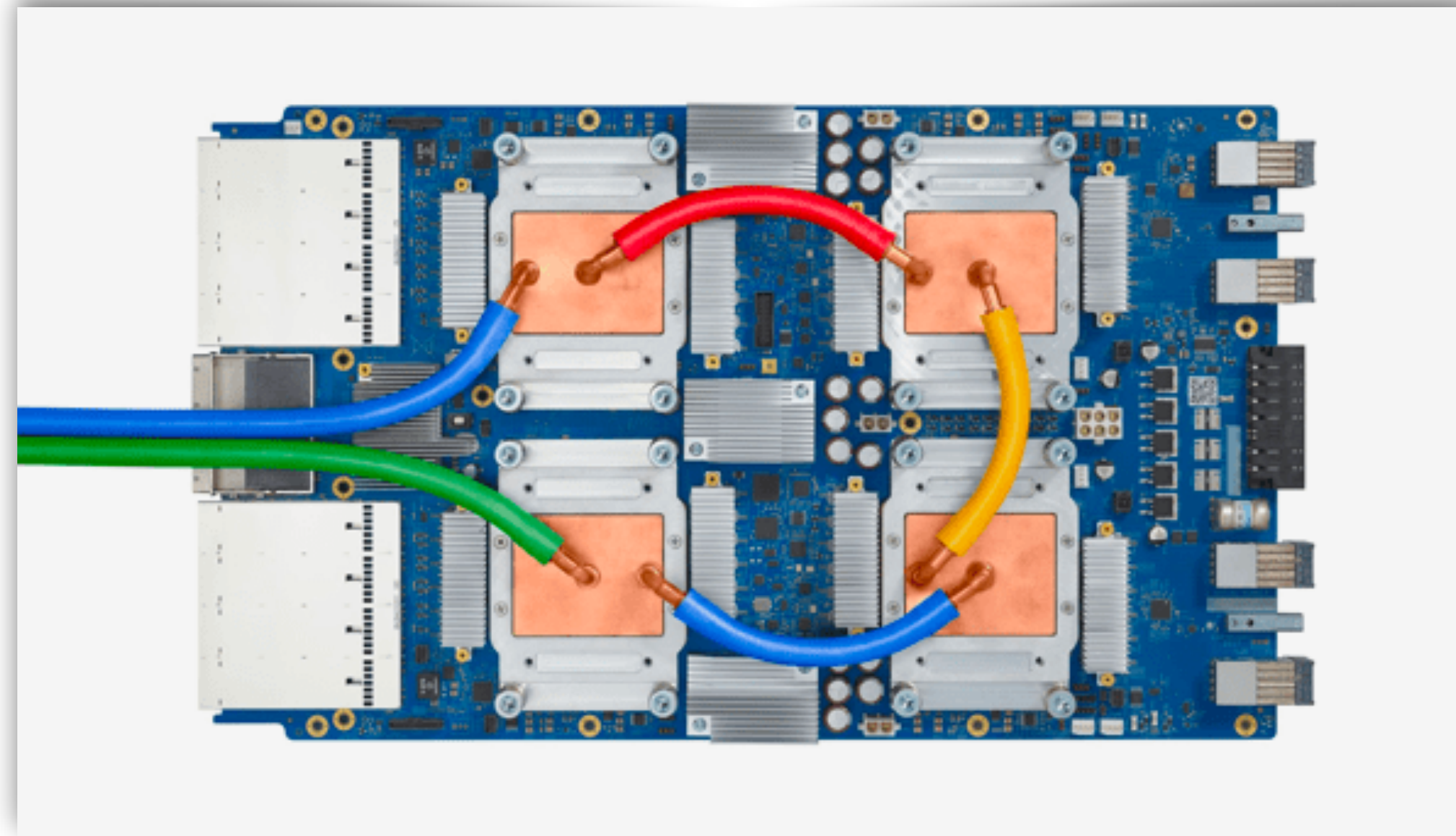






Should be fast and power-efficient

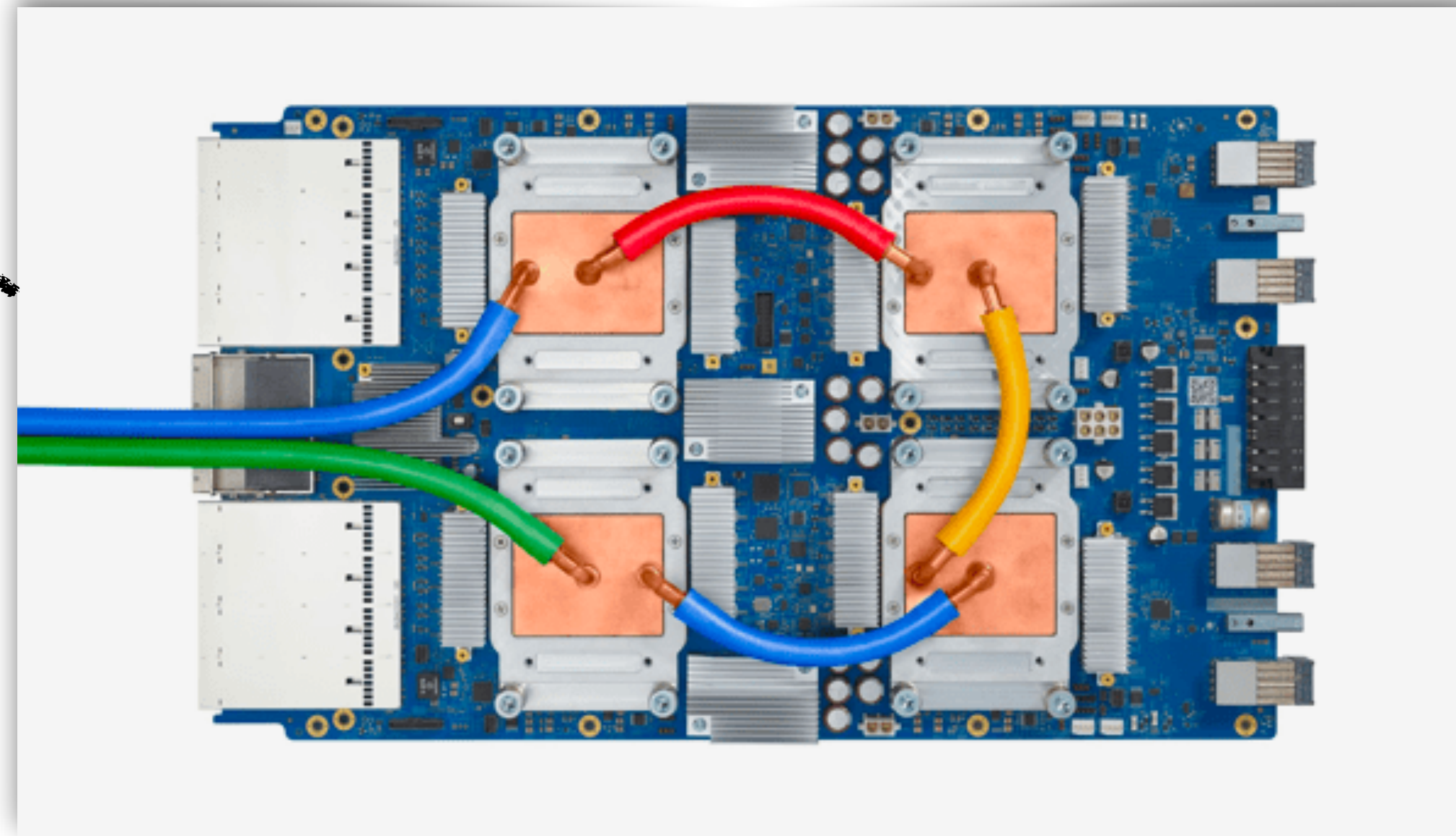




Should be fast and power-efficient

Needs small weights and activations
to maximize usage of chip area

Only needs to represent a specific range of values:
Weights and activations cluster (e.g. around $[-1, 1]$)

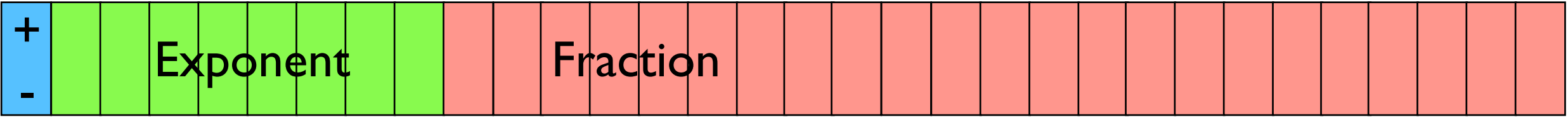


Should be fast and power-efficient

Needs small weights and activations
to maximize usage of chip area

bfloat16

Single precision IEEE float (32 bit)

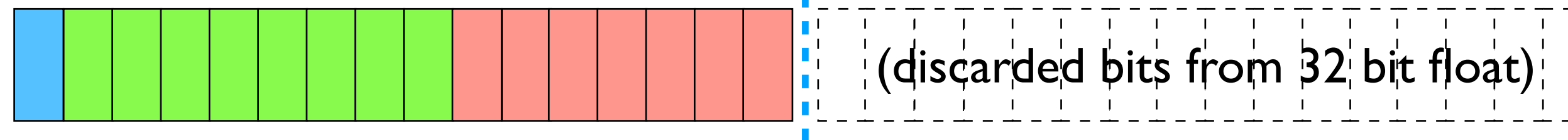


bfloat16

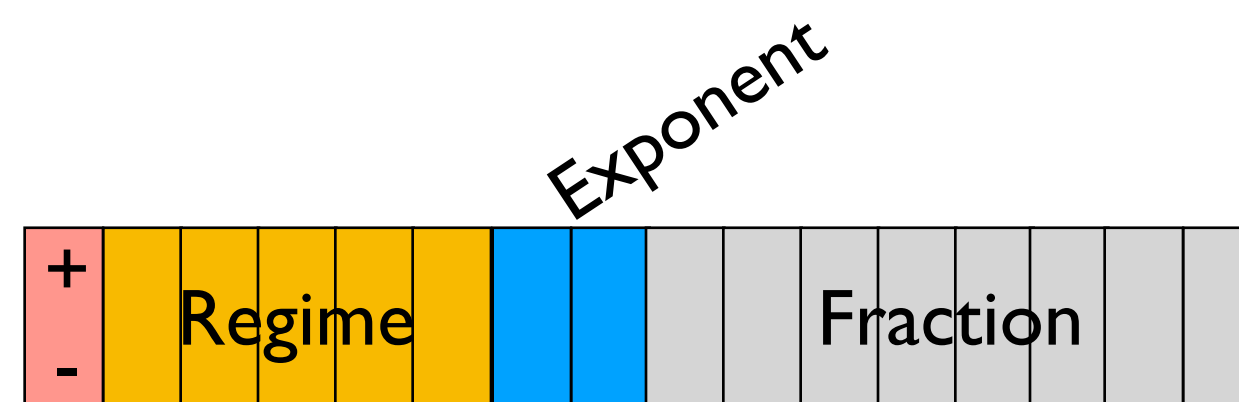
Single precision IEEE float (32 bit)



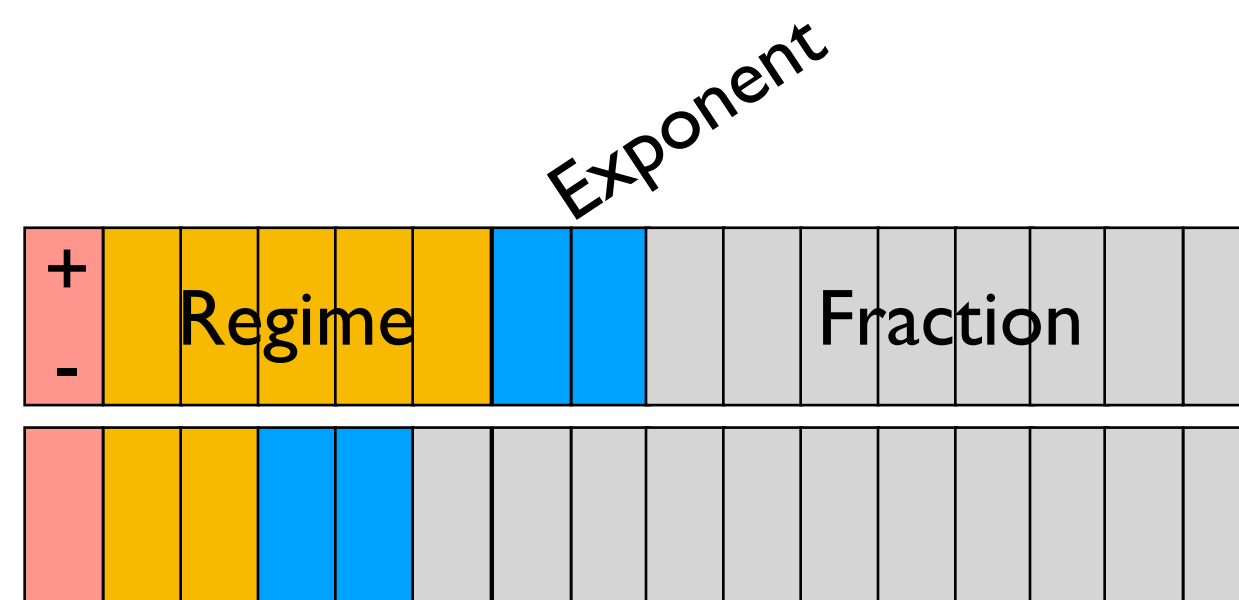
bfloat16



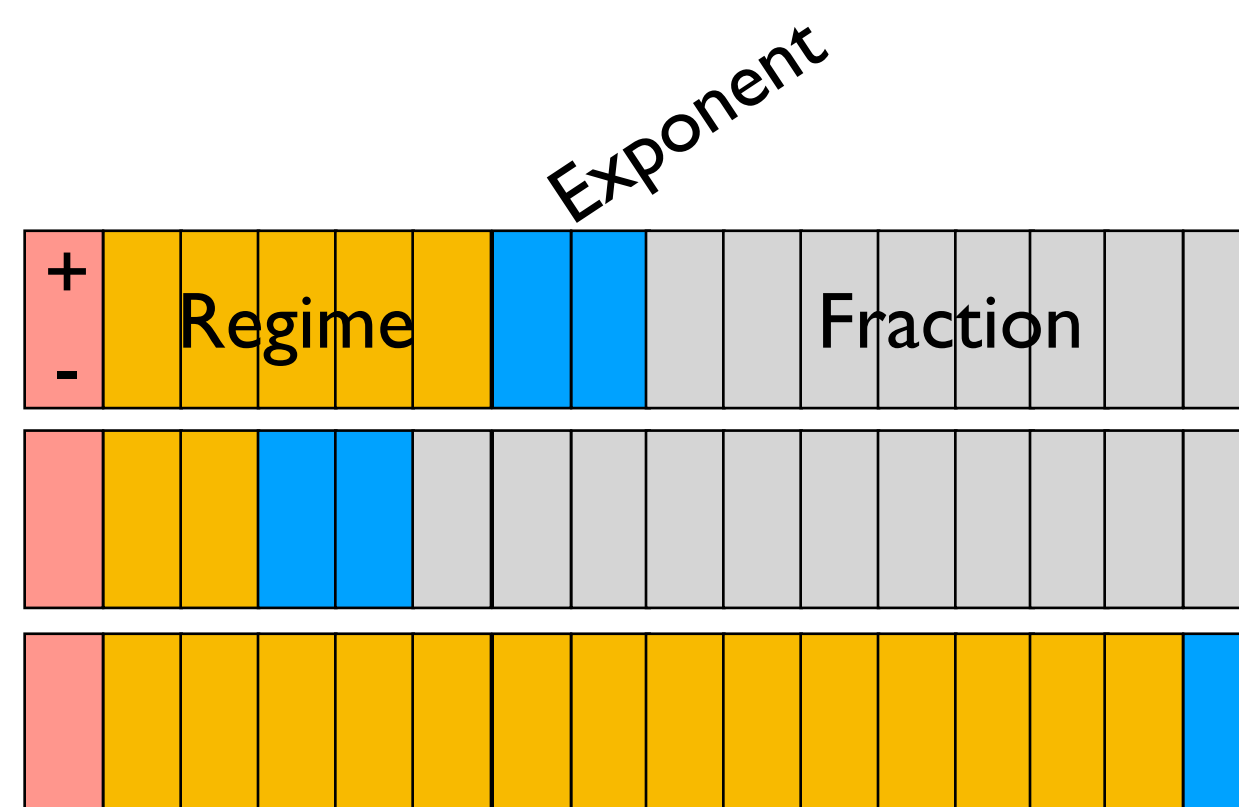
Posits



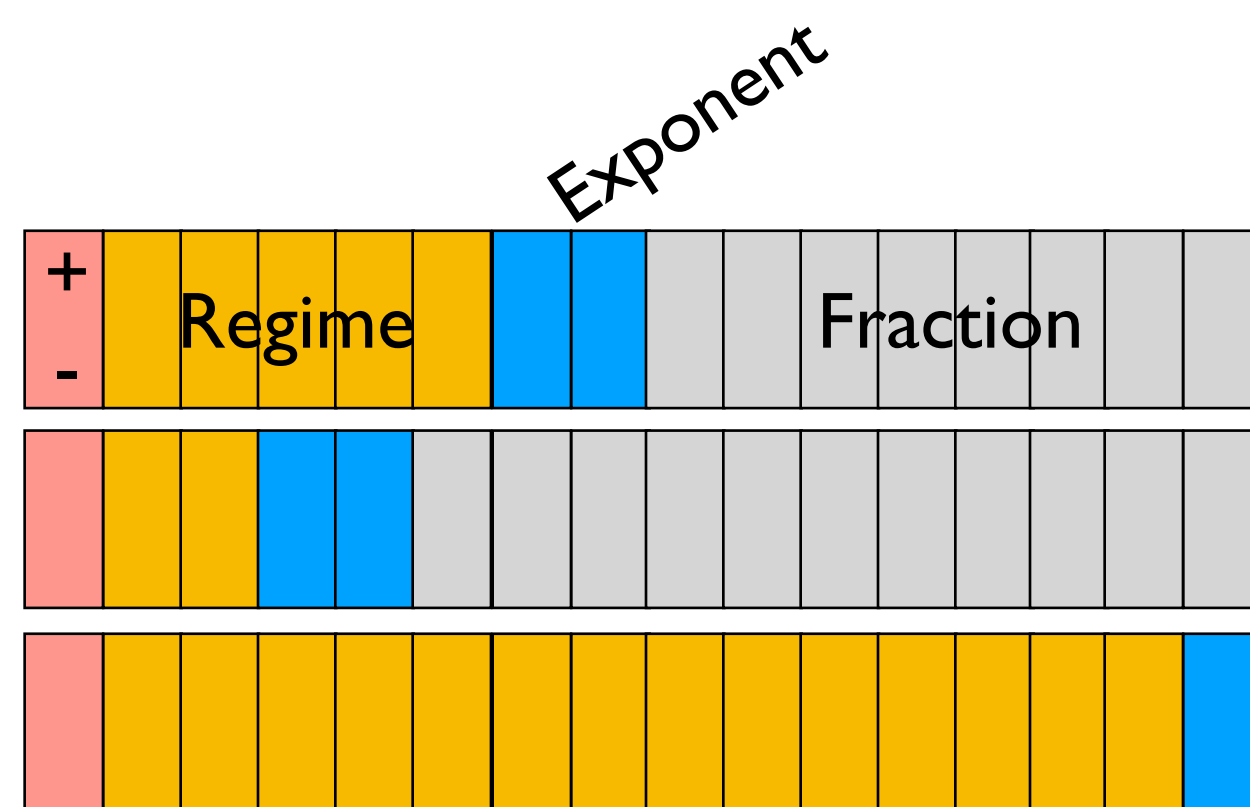
Posits



Posits

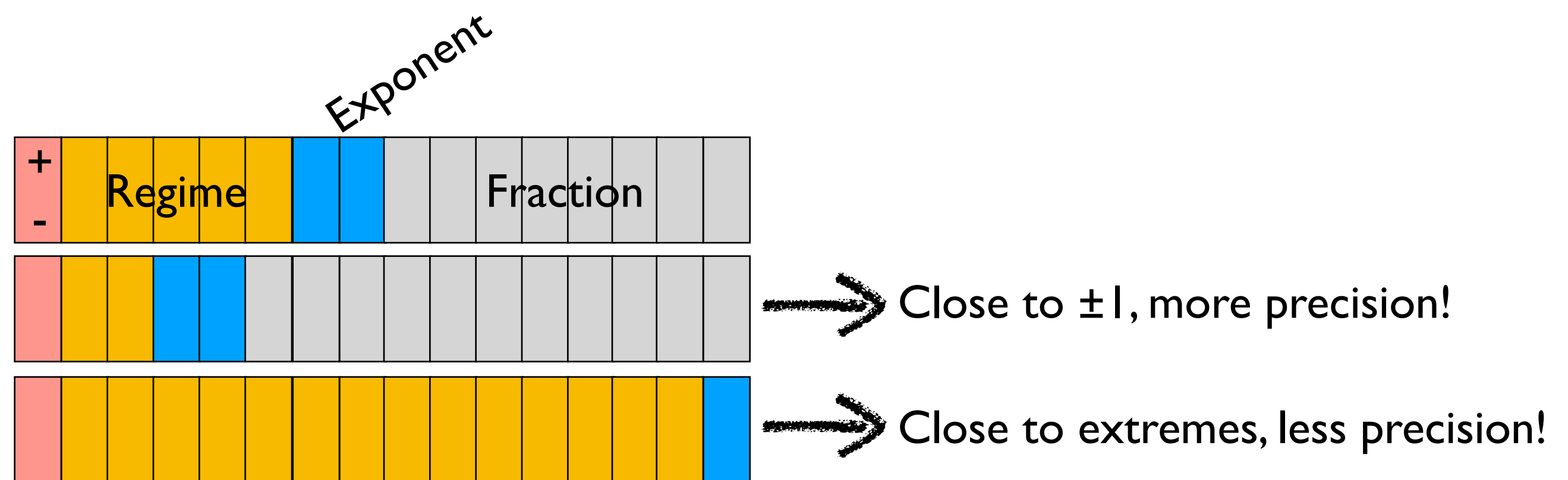


Posits



→ Close to ± 1 , more precision!

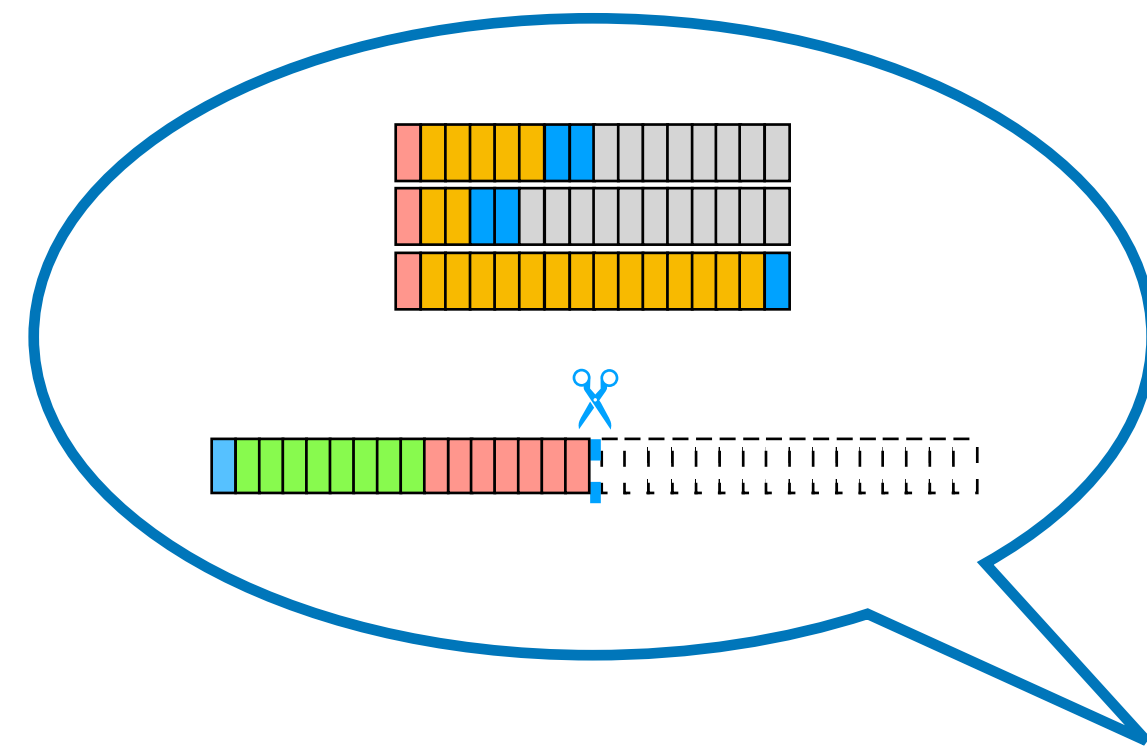
Posits

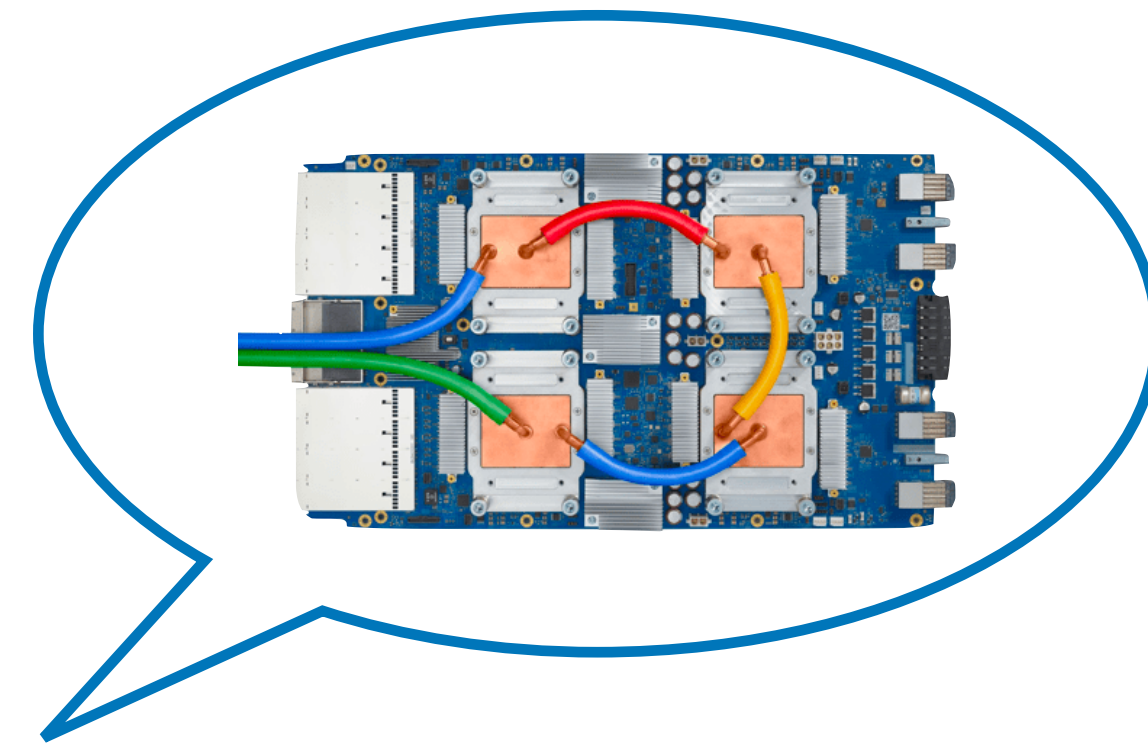
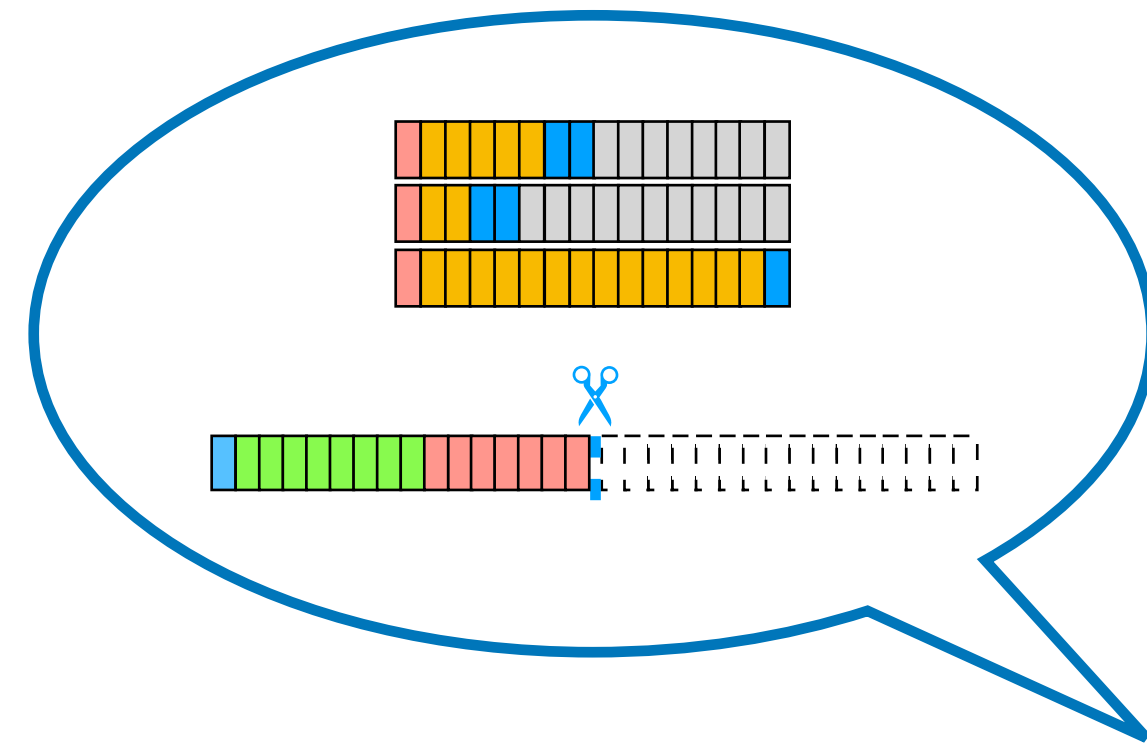


Deep learning needs hardware specialization...

Deep learning needs hardware specialization...
...and hardware specialization needs new datatypes!







picopicodevil / bfloat16

Watch

1

Star

0

Fork

0

<> Code

Issues 0

Pull requests 0

Actions

Projects 0

Wiki

Security

Insights

No description, website, or topics provided.

1 commit

1 branch

N-Dekker / biovault_bfloat16

Watch

1

Star

0

Fork

1

<> Code

Issues 0

Pull requests 0

Actions

Projects 0

Wiki

Security

Insights

A bfloat16 implementation for BioVault projects

1 contributor

Apache-2.0

libcg / bfp

Watch

31

Star

218

Fork

21

<> Code

Issues 2

Pull requests 0

Actions

Projects 0

Wiki

Security

Insights

Beyond Floating Point - Posit C/C++ implementation

posit

gustafson

ieee754

unum

122 commits

stillwater-sc / universal

Watch

19

Unstar

125

Fork

16

<> Code

Issues 8

Pull requests 0

Actions

Projects 3

Wiki

Security

Insights

tors

MIT

cjdelisle / libposit

Watch

1

Unstar

7

Fork

1

<> Code

Issues 0

Pull requests 1

Actions

Projects 0

Wiki

Security

Insights

A library for working with the posit number type.

46 commits

1 branch

0 packages

0 releases

1 contributor

View license

Branch: master

New pull request

Create new file

Upload files

Find file

Clone or download

Types

libposit defines the following types:

- `posit8_t` 8 bit posit
 - `posit8x2_t` pair of 8 bit posits, result of `*_exact` functions.
- `posit16_t` 16 bit posit
 - `posit16x2_t` pair of 16 bit posits, result of `*_exact` functions.
- `posit32_t` 32 bit posit
 - `posit32x2_t` pair of 32 bit posits, result of `*_exact` functions.
- `posit64_t` 64 bit posit
 - `posit64x2_t` pair of 64 bit posits, result of `*_exact` functions.
- `posit128_t` 128 bit posit, result of posit64 `*_promote` functions. Only functions which work on a `posit128_t` are `posit128_iexp()` and `posit128_fract()`.

Types

libposit defines the

- `posit8_t` 8 bits
 - `posit8x`
- `posit16_t` 16 bits
 - `posit16x`
- `posit32_t` 32 bits
 - `posit32x`
- `posit64_t` 64 bits
 - `posit64x`
- `posit128_t` 128 bits
 - `posit128_ie`

Operations

For each of these types, libposit provides the following functions:

- `posit<X>_t posit<X>_add(posit<X>_t x, posit<X>_t y)` : Add two posits together, output a rounded result.
- `posit<X>_t posit<X>_sub(posit<X>_t x, posit<X>_t y)` : Subtract one posit from another, output a rounded result.
- `posit<X>x2_t posit<X>_add_exact(posit<X>_t x, posit<X>_t y)` : Output 2 results, `sum` and `remainder` which when added will result in the exact same value as the inputs `x` and `y` but where the absolute value of `remainder` is as small as possible.
- `posit<X>x2_t posit<X>_sub_exact(posit<X>_t x, posit<X>_t y)` : Exactly the same as `posit<X>_add_exact(x, y)` but with the sign of `y` flipped.
- `posit<X>_t posit<X>_mul(posit<X>_t x, posit<X>_t y)` : Multiply two posits, round the result to nearest
- `posit<X>_t posit<X>_div(posit<X>_t x, posit<X>_t y)` : Divide two posits, round the result to nearest

Types

libposit defines the

- `posit8_t` 8 bits
 - `posit8x`
- `posit16_t` 16 bits
 - `posit16x`
- `posit32_t` 32 bits
 - `posit32x`
- `posit64_t` 64 bits
 - `posit64x`
- `posit128_t` 128 bits
 - `posit128_ie`

Operations

For each of these types, libposit provides the following functions:

- `posit<X>_t posit<X>_add(posit<X>_t x, posit<X>_t y)` : Add two posits together, output a rounded result.
- `posit<X>_t posit<X>_sub(posit<X>_t x, posit<X>_t y)` : Subtract one posit from another, output a rounded result.
- `posit<X>x2_t posit<X>_add_exact(posit<X>_t x, posit<X>_t y)` !: Output 2 results, `sum` and `remainder` which when added will result in the exact same value as the inputs `x` and `y` but where the absolute value of `remainder` is as small as possible.
- `posit<X>x2_t posit<X>_sub_exact(posit<X>_t x, posit<X>_t y)` : Exactly the same as `posit<X>_add_exact(x, y)` but with the sign of `y` flipped.
- `posit<X>_t posit<X>_mul(posit<X>_t x, posit<X>_t y)` : Multiply two posits, round the result to nearest
- `posit<X>_t posit<X>_div(posit<X>_t x, posit<X>_t y)` : Divide two posits, round the result to nearest

Types

libposit defines the

- `posit8_t` 8 bits
 - `posit8x`
- `posit16_t` 16 bits
 - `posit16x`
- `posit32_t` 32 bits
 - `posit32x`
- `posit64_t` 64 bits
 - `posit64x`
- `posit128_t` 128 bits
 - `posit128_ie`

Operations

For each of these types, libposit provides the following functions:

- `posit<X>_t posit<X>_add(posit<X>_t x, posit<X>_t y)` : Add two posits together, output a rounded result.
- `posit<X>_t posit<X>_sub(posit<X>_t x, posit<X>_t y)` : Subtract one posit from another, output a rounded result.
- `posit<X>x2_t posit<X>_add_exact(posit<X>_t x, posit<X>_t y)` !: Output 2 results, `sum` and `remainder` which when added will result in the exact same value as the inputs `x` and `y` but where the absolute value of `remainder` is as small as possible.
- `posit<X>x2_t posit<X>_sub_exact(posit<X>_t x, posit<X>_t y)` !: Exactly the same as `posit<X>_add_exact(x, y)` but with the sign of `y` flipped.
- `posit<X>_t posit<X>_mul(posit<X>_t x, posit<X>_t y)` : Multiply two posits, round the result to nearest
- `posit<X>_t posit<X>_div(posit<X>_t x, posit<X>_t y)` : Divide two posits, round the result to nearest

Types

libposit defines the

- `posit8_t` 8 bits
 - `posit8x`
- `posit16_t` 16 bits
 - `posit16x`
- `posit32_t` 32 bits
 - `posit32x`
- `posit64_t` 64 bits
 - `posit64x`
- `posit128_t` 128 bits
 - `posit128_ie`

Operations

For each of these types, libposit provides the following functions:

- `posit<X>_t posit<X>_add(posit<X>_t x, posit<X>_t y)` ✓ Add two posits together, output a rounded result.
- `posit<X>_t posit<X>_sub(posit<X>_t x, posit<X>_t y)` : Subtract one posit from another, output a rounded result.
- `posit<X>x2_t posit<X>_add_exact(posit<X>_t x, posit<X>_t y)` !: Output 2 results, `sum` and `remainder` which when added will result in the exact same value as the inputs `x` and `y` but where the absolute value of `remainder` is as small as possible.
- `posit<X>x2_t posit<X>_sub_exact(posit<X>_t x, posit<X>_t y)` !: Exactly the same as `posit<X>_add_exact(x, y)` but with the sign of `y` flipped.
- `posit<X>_t posit<X>_mul(posit<X>_t x, posit<X>_t y)` : Multiply two posits, round the result to nearest
- `posit<X>_t posit<X>_div(posit<X>_t x, posit<X>_t y)` : Divide two posits, round the result to nearest

Types

libposit defines the

- `posit8_t` 8 bits
 - `posit8x`
- `posit16_t` 16 bits
 - `posit16x`
- `posit32_t` 32 bits
 - `posit32x`
- `posit64_t` 64 bits
 - `posit64x`
- `posit128_t` 128 bits
 - `posit128_ie`

Operations

For each of these types, libposit provides the following functions:

- `posit<X>_t posit<X>_add(posit<X>_t x, posit<X>_t y)` ✓ Add two posits together, output a rounded result.
- `posit<X>_t posit<X>_sub(posit<X>_t x, posit<X>_t y)` ✓ Subtract one posit from another, output a rounded result.
- `posit<X>x2_t posit<X>_add_exact(posit<X>_t x, posit<X>_t y)` !: Output 2 results, `sum` and `remainder` which when added will result in the exact same value as the inputs `x` and `y` but where the absolute value of `remainder` is as small as possible.
- `posit<X>x2_t posit<X>_sub_exact(posit<X>_t x, posit<X>_t y)` !: Exactly the same as `posit<X>_add_exact(x, y)` but with the sign of `y` flipped.
- `posit<X>_t posit<X>_mul(posit<X>_t x, posit<X>_t y)` : Multiply two posits, round the result to nearest
- `posit<X>_t posit<X>_div(posit<X>_t x, posit<X>_t y)` : Divide two posits, round the result to nearest

Types

libposit defines the

- `posit8_t` 8 bits
 - `posit8x`
- `posit16_t` 16 bits
 - `posit16x`
- `posit32_t` 32 bits
 - `posit32x`
- `posit64_t` 64 bits
 - `posit64x`
- `posit128_t` 128 bits
 - `posit128_ie`

Operations

For each of these types, libposit provides the following functions:

- `posit<X>_t posit<X>_add(posit<X>_t x, posit<X>_t y)` ✓ Add two posits together, output a rounded result.
- `posit<X>_t posit<X>_sub(posit<X>_t x, posit<X>_t y)` ✓ Subtract one posit from another, output a rounded result.
- `posit<X>x2_t posit<X>_add_exact(posit<X>_t x, posit<X>_t y)` !: Output 2 results, `sum` and `remainder` which when added will result in the exact same value as the inputs `x` and `y` but where the absolute value of `remainder` is as small as possible.
- `posit<X>x2_t posit<X>_sub_exact(posit<X>_t x, posit<X>_t y)` !: Exactly the same as `posit<X>_add_exact(x, y)` but with the sign of `y` flipped.
- `posit<X>_t posit<X>_mul(posit<X>_t x, posit<X>_t y)` ✓ Multiply two posits, round the result to nearest
- `posit<X>_t posit<X>_div(posit<X>_t x, posit<X>_t y)` : Divide two posits, round the result to nearest

Types

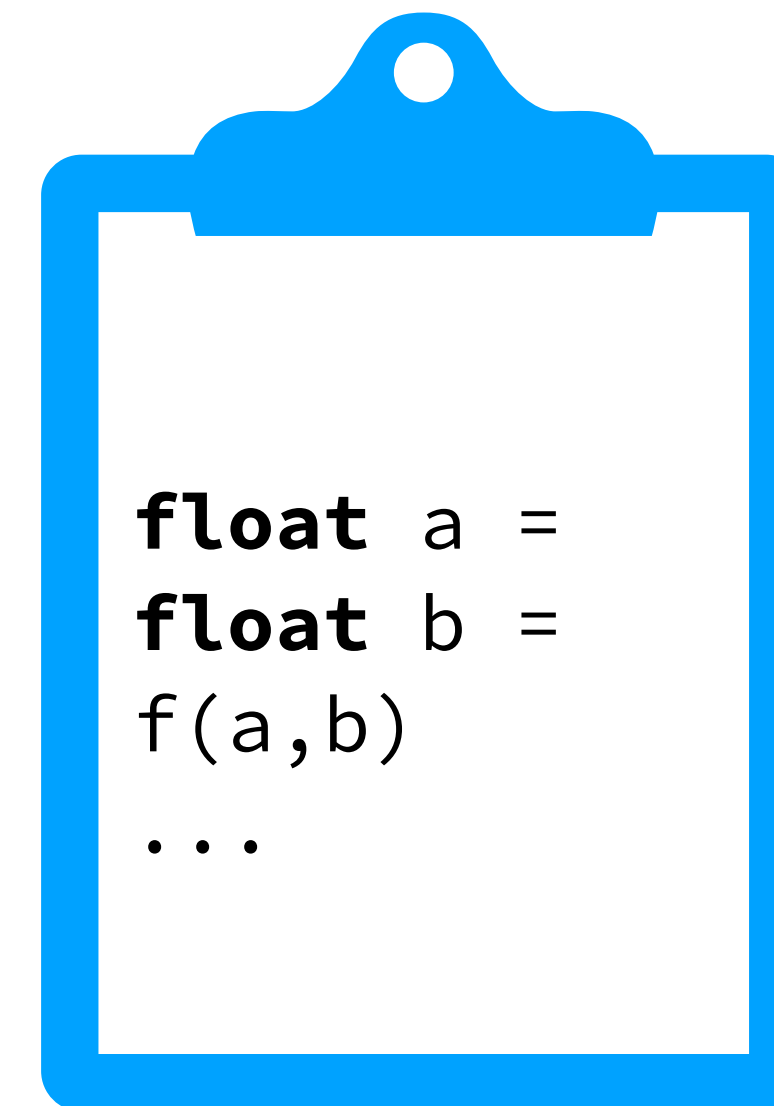
libposit defines the

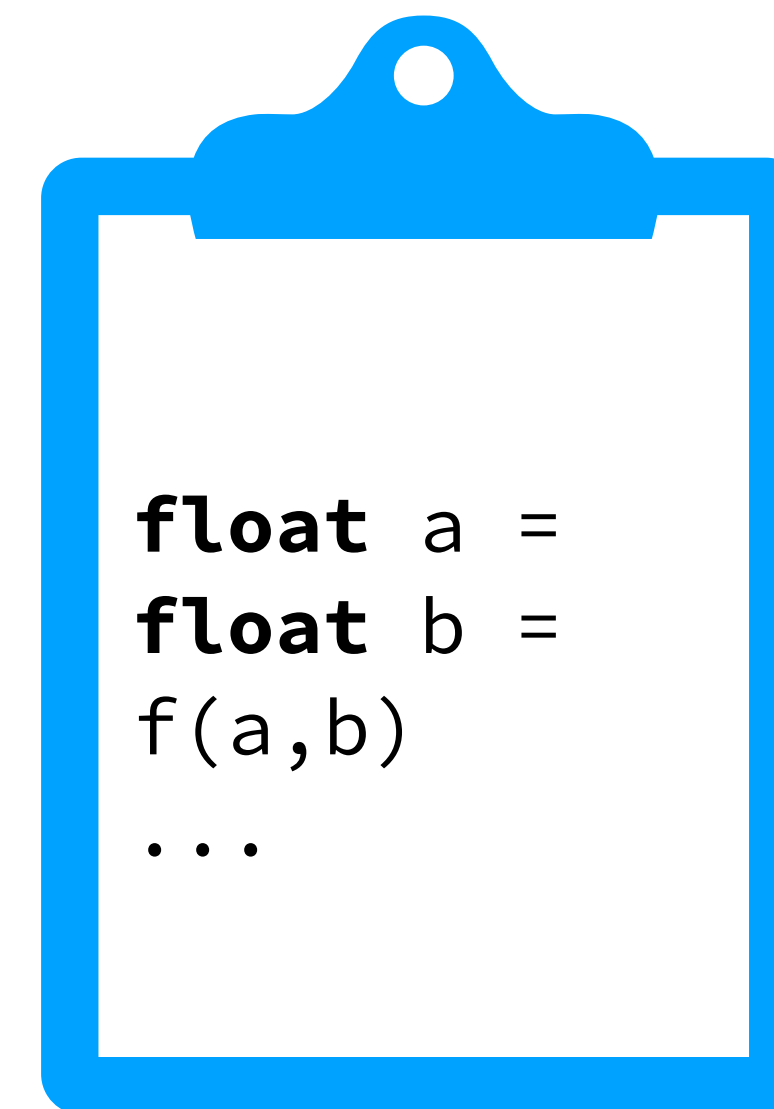
- `posit8_t` 8 bits
 - `posit8x`
- `posit16_t` 16 bits
 - `posit16x`
- `posit32_t` 32 bits
 - `posit32x`
- `posit64_t` 64 bits
 - `posit64x`
- `posit128_t` 128 bits
 - `posit128_ie`

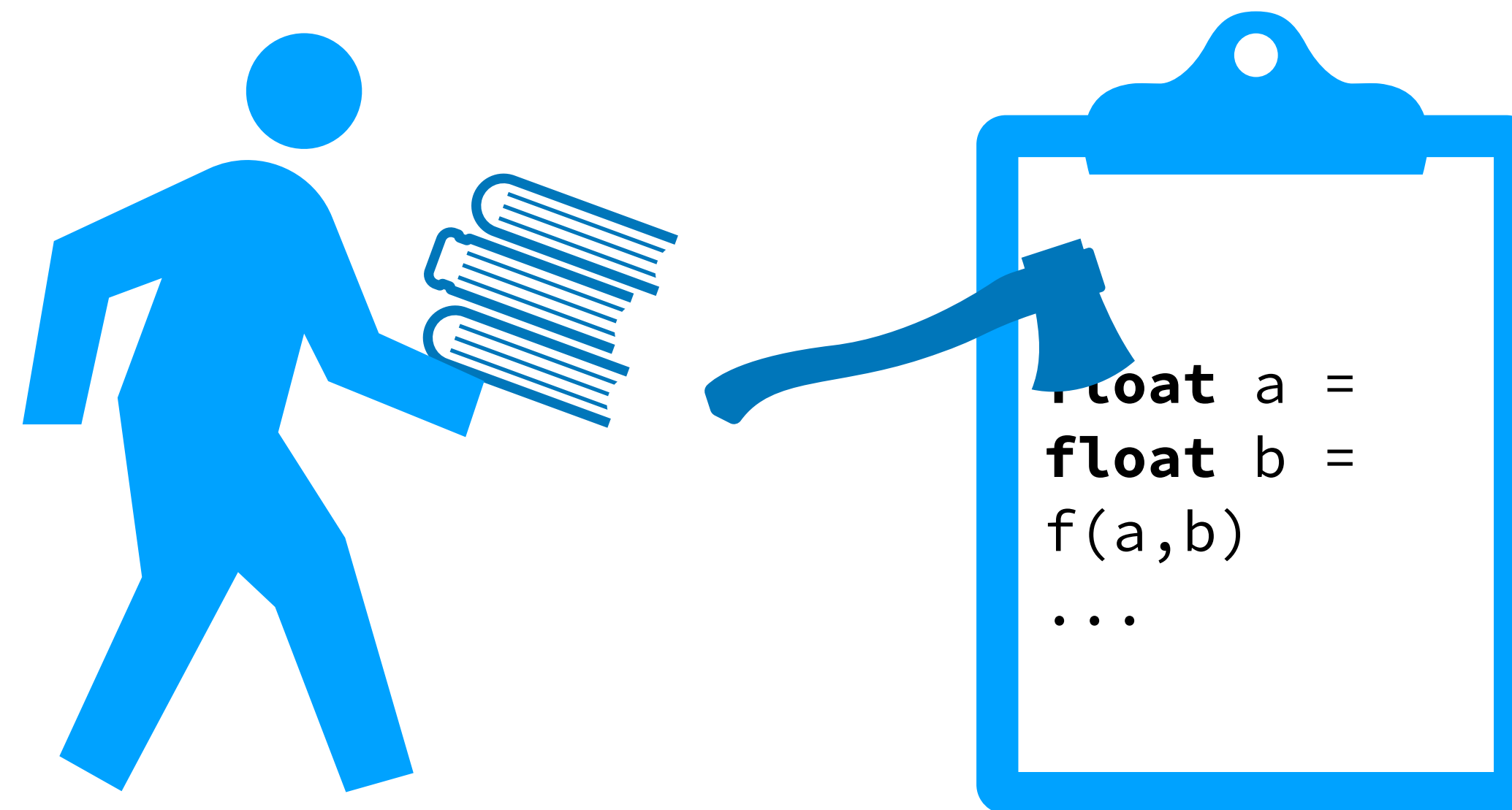
Operations

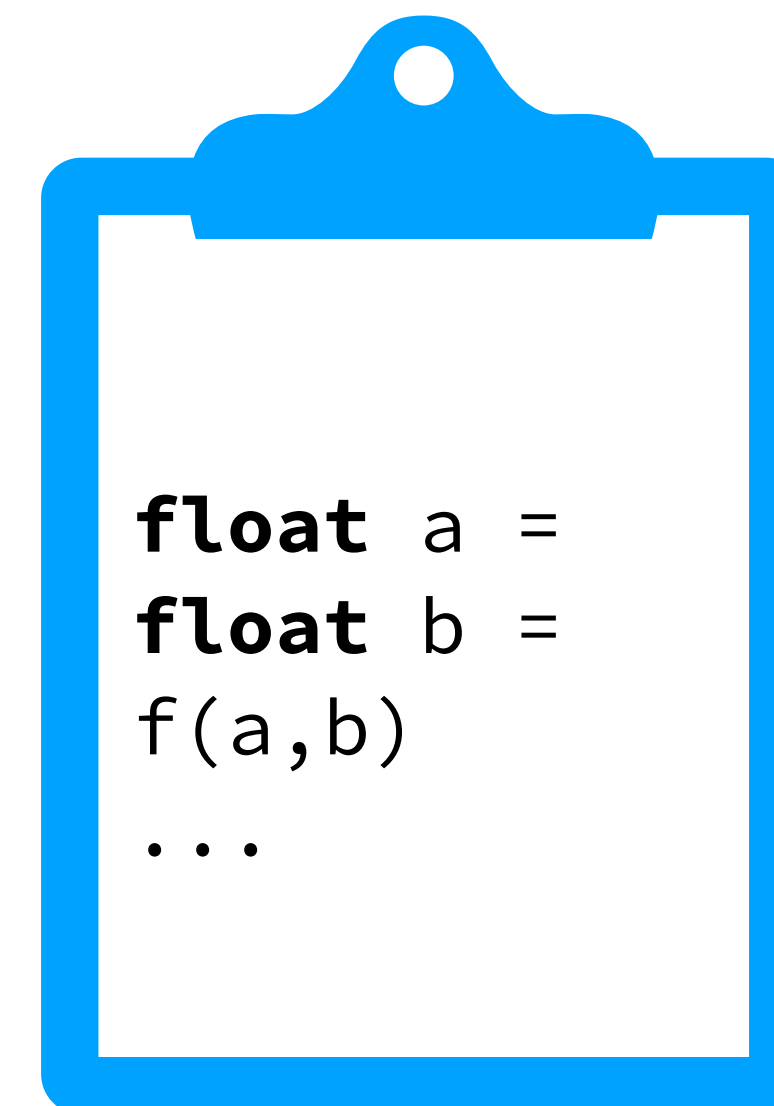
For each of these types, libposit provides the following functions:

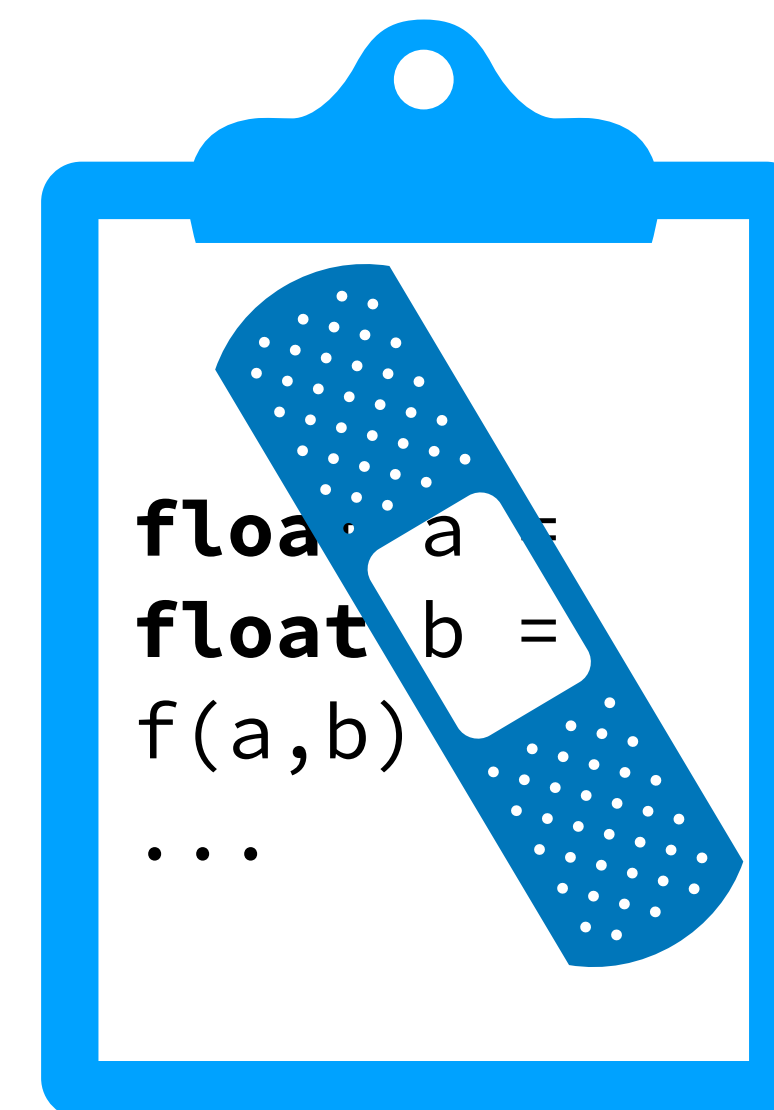
- `posit<X>_t posit<X>_add(posit<X>_t x, posit<X>_t y)` ✓ Add two posits together, output a rounded result.
- `posit<X>_t posit<X>_sub(posit<X>_t x, posit<X>_t y)` ✓ Subtract one posit from another, output a rounded result.
- `posit<X>x2_t posit<X>_add_exact(posit<X>_t x, posit<X>_t y)` !: Output 2 results, `sum` and `remainder` which when added will result in the exact same value as the inputs `x` and `y` but where the absolute value of `remainder` is as small as possible.
- `posit<X>x2_t posit<X>_sub_exact(posit<X>_t x, posit<X>_t y)` !: Exactly the same as `posit<X>_add_exact(x, y)` but with the sign of `y` flipped.
- `posit<X>_t posit<X>_mul(posit<X>_t x, posit<X>_t y)` ✓ Multiply two posits, round the result to nearest
- `posit<X>_t posit<X>_div(posit<X>_t x, posit<X>_t y)` ✓ Divide two posits, round the result to nearest

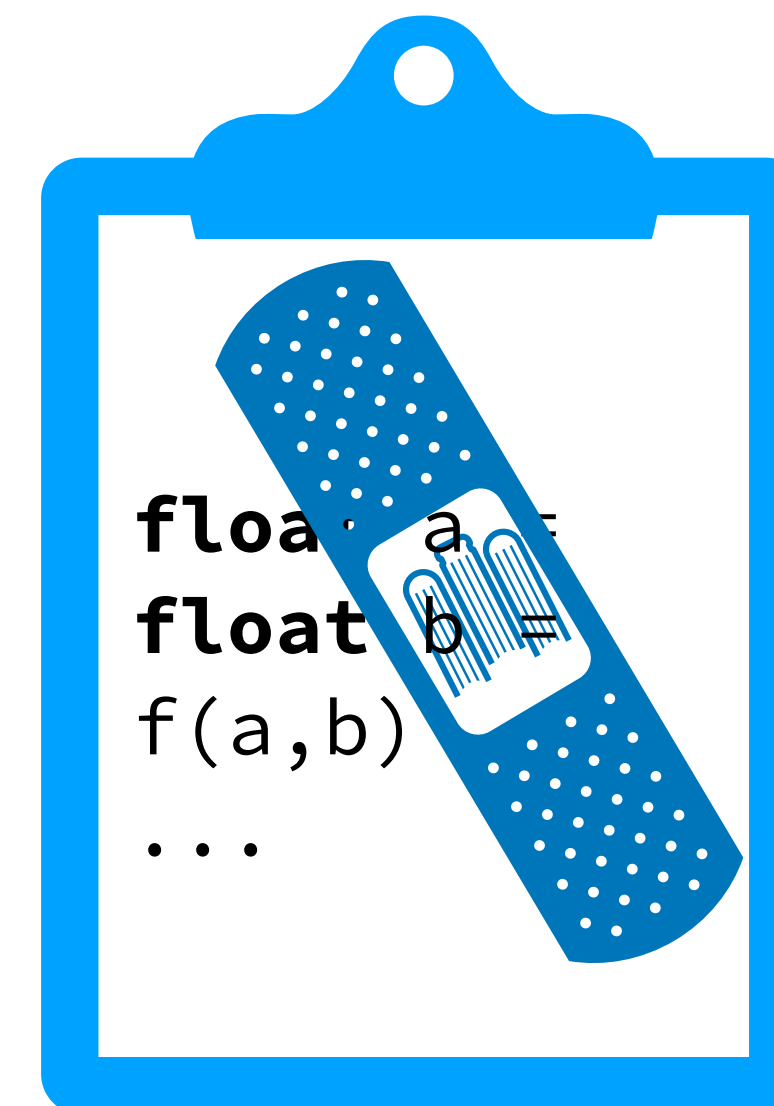


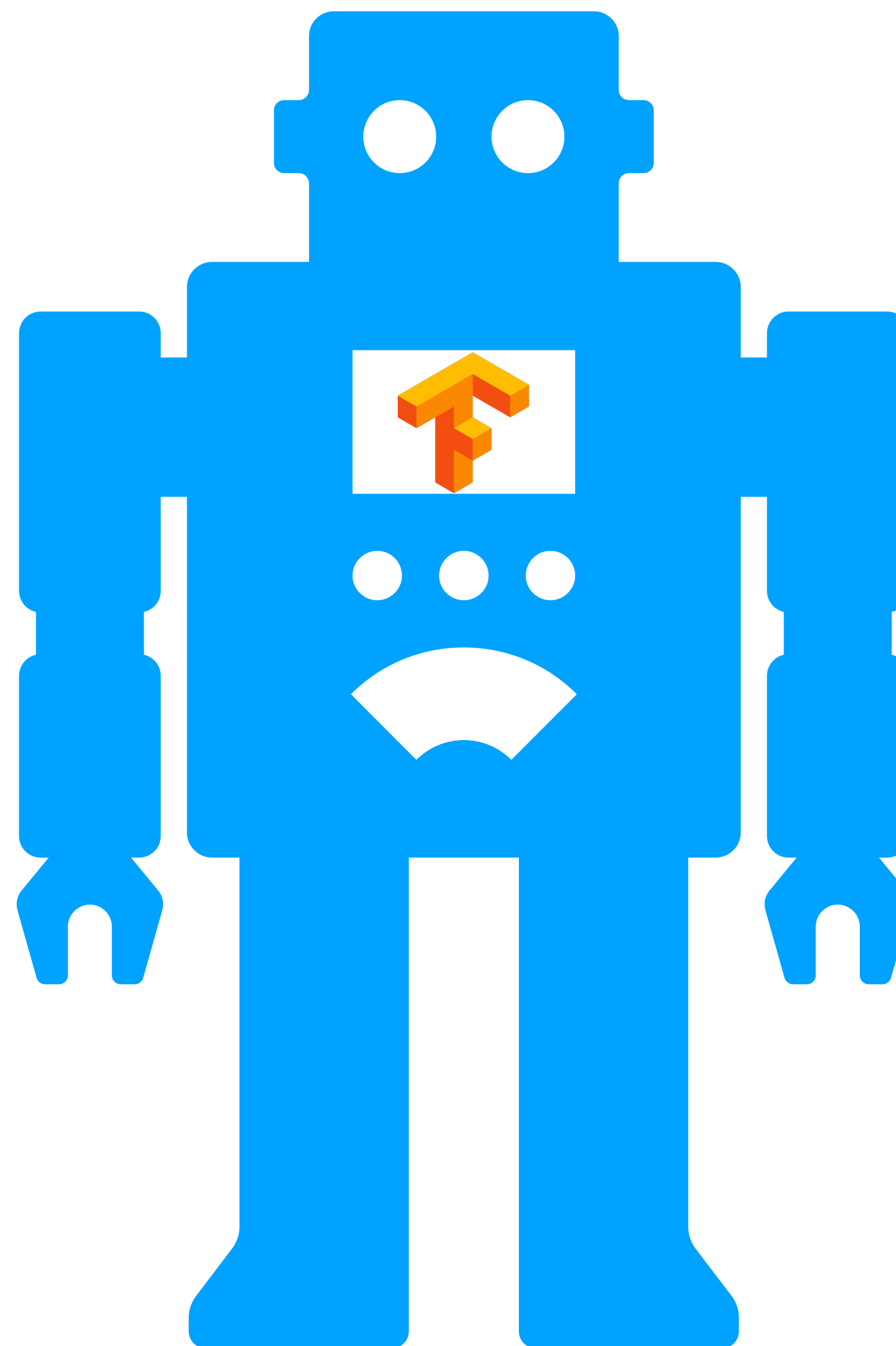




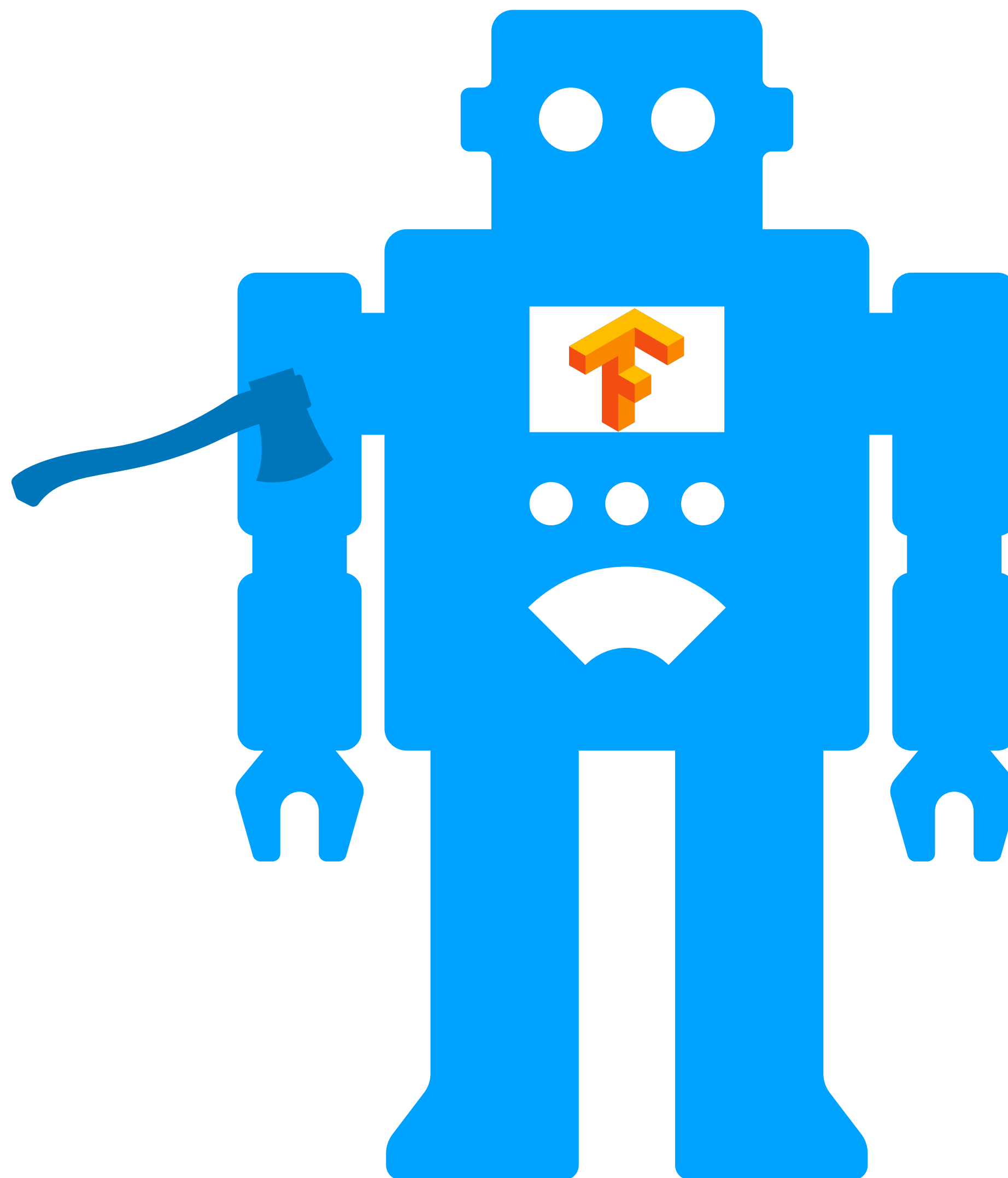




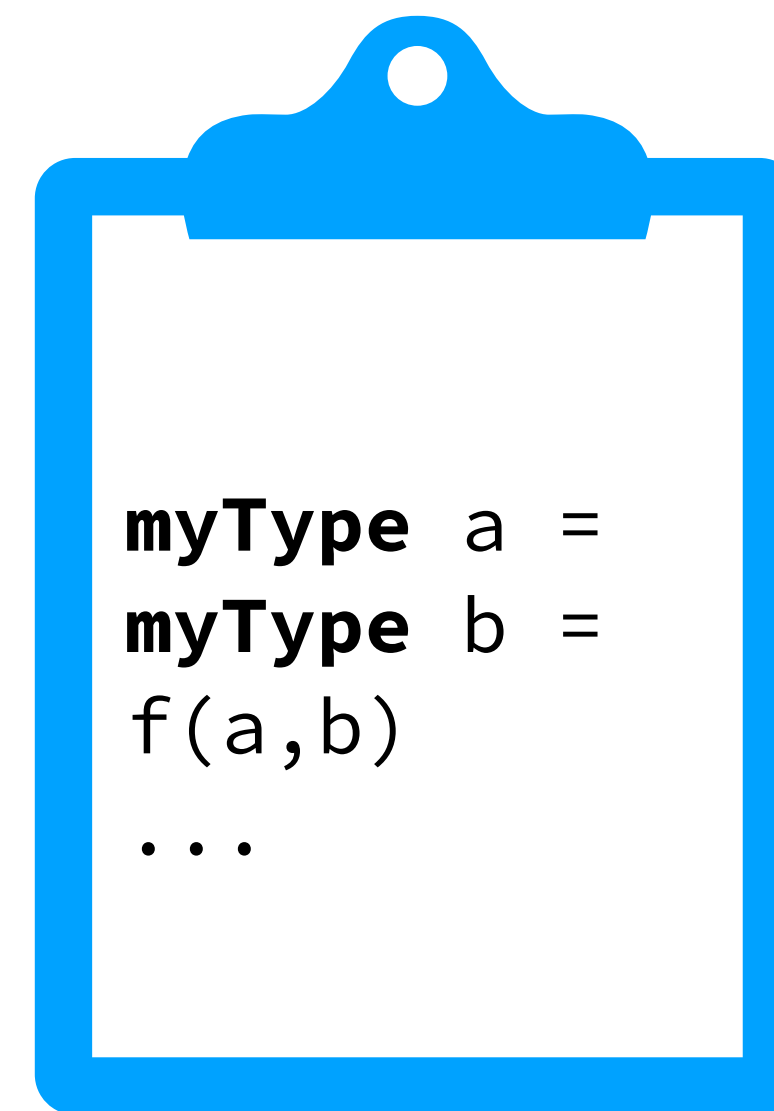
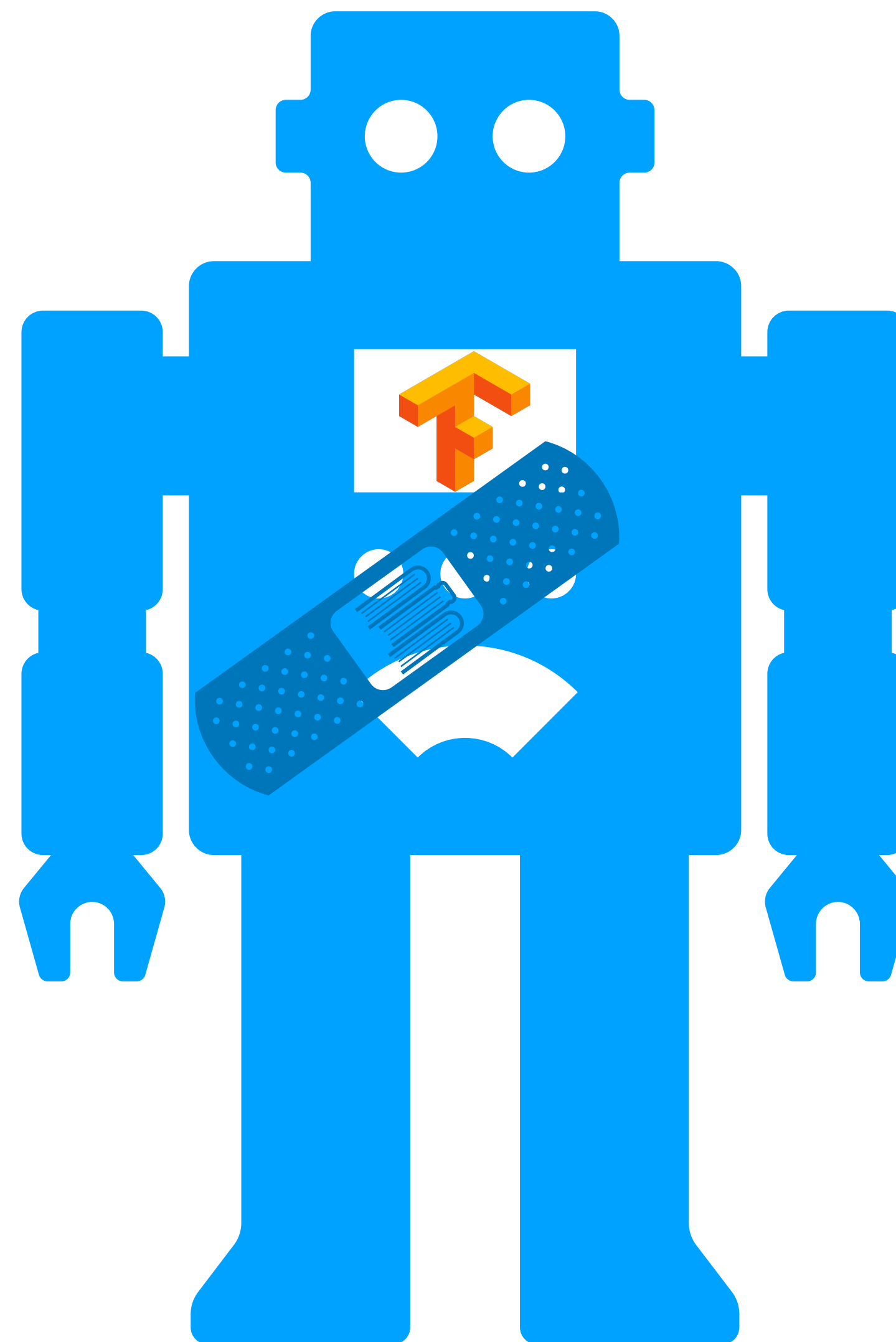


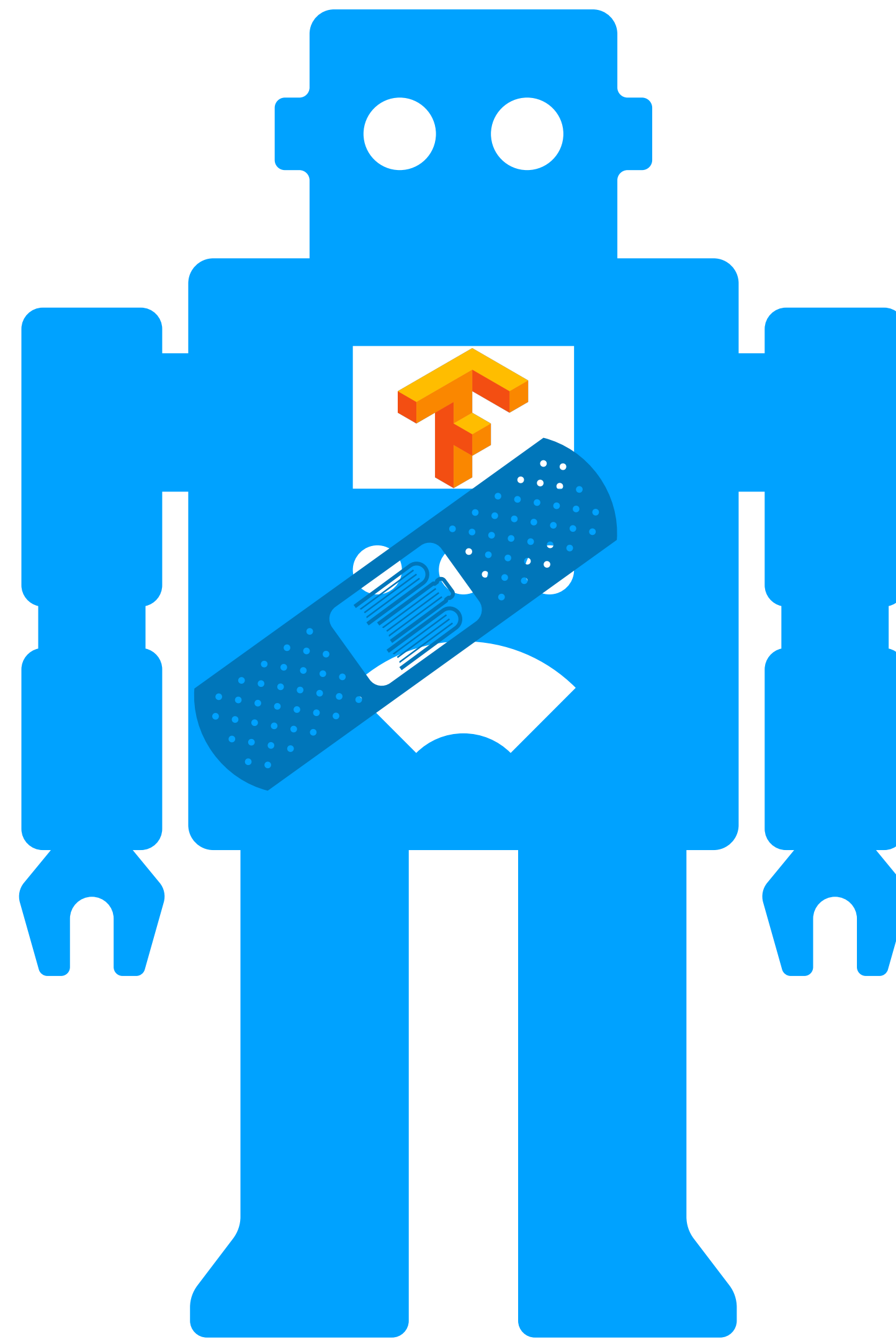


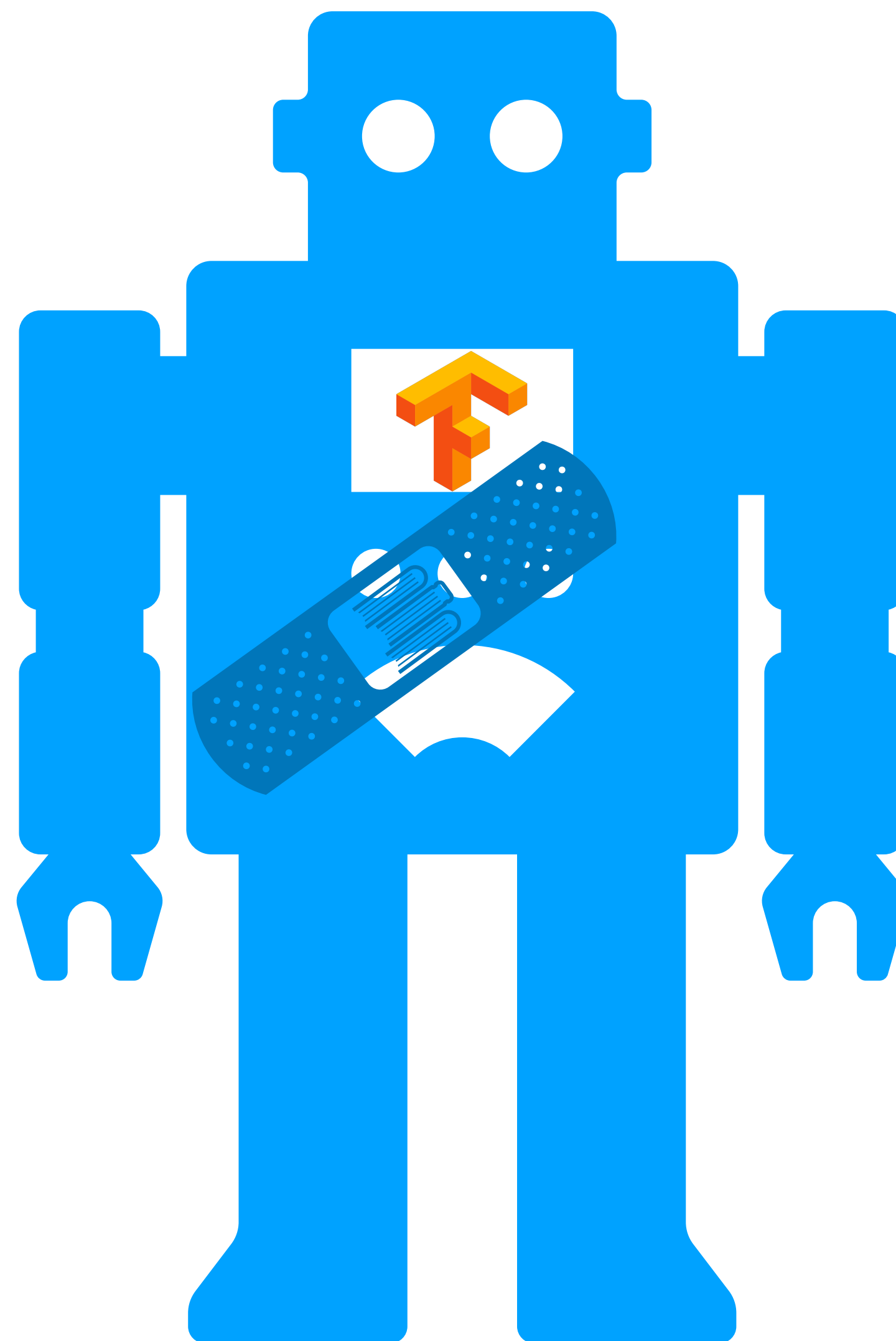
```
myType a =  
myType b =  
f(a,b)  
...
```



```
myType a =  
myType b =  
f(a,b)  
...
```











An Example in the Wild!


An Example in the Wild!


 **xman / tensorflow**
forked from [tensorflow/tensorflow](#)


 Watch 3

 Star 5


 Fork 75,570


 Code


 Issues 0

 Pull requests 0


 Projects 0




 Wiki








 Security






 Insights






An Example in the Wild!



 **xman / tensorflow**
forked from [tensorflow/tensorflow](#)



 Watch 3  Star 5  Fork 75,570

 Code  Issues 0  Pull requests 0  Projects 0  Wiki  Security  Insights


 36,645 commits  14 branches  8 releases  1,481 contributors  Apache-2.0




Branch: **posit**   New pull request  Find File  Clone or download 








This branch is **237 commits ahead**, 22081 commits behind tensorflow:master.  Pull request  Compare






 **xman** posit: update version.  Latest commit 47fa4c7 on Sep 20, 2018






An Example in the Wild!



 **xman / tensorflow**
forked from [tensorflow/tensorflow](#)



 Watch 3  Star 5  Fork 75,570




 Code  Issues 0  Pull requests 0  Projects 0  Wiki  Security  Insights

 36,645 commits  14 branches  8 releases  1,481 contributors  Apache-2.0


Branch: **posit**   New pull request  Find File  Clone or download 




This branch is **237 commits ahead**, 22081 commits behind tensorflow:master.  Pull request  Compare








 **xman** posit: update version.  Latest commit 47fa4c7 on Sep 20, 2018






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






An Example in the Wild!



 **xman / tensorflow**
forked from [tensorflow/tensorflow](#)



 Watch 3  Star 5  Fork 75,570




 Code  Issues 0  Pull requests 0  Projects 0  Wiki  Security  Insights

 36,645 commits  14 branches  8 releases  1,481 contributors  Apache-2.0

Branch: **posit**   New pull request  Find File  Clone or download 

This branch is **237 commits ahead**, 22081 commits behind tensorflow:master.  Pull request  Compare

 **xman** posit: update version.  Latest commit 47fa4c7 on Sep 20, 2018

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

The computational demands of deep learning require new datatypes...

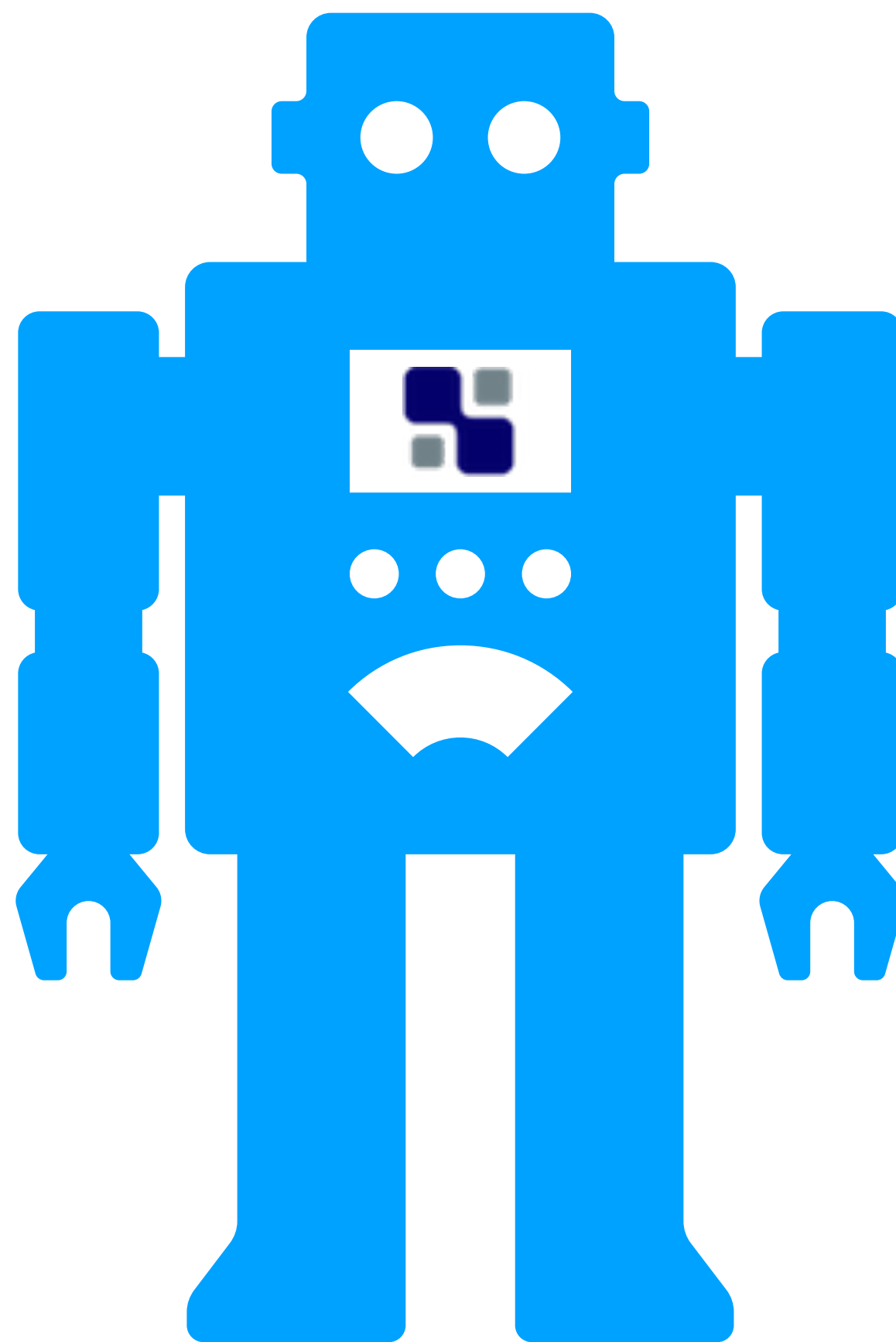
The computational demands of deep learning require new datatypes...

...but datatype research is difficult!

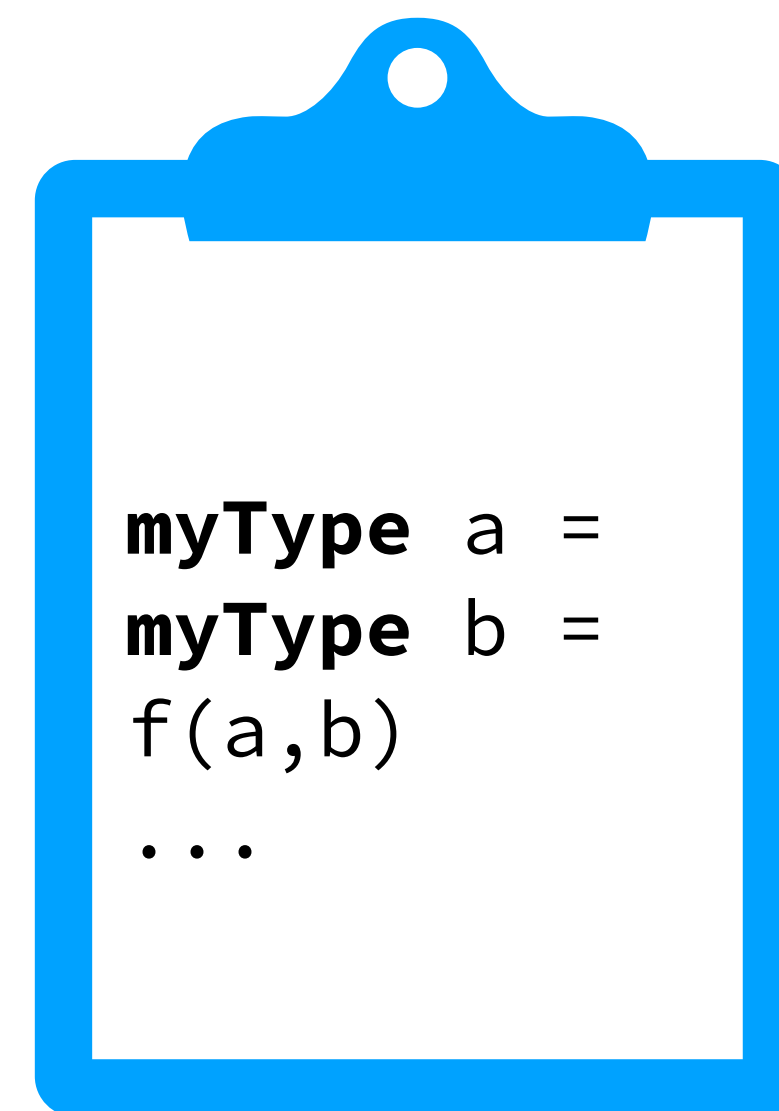
The computational demands of deep learning require new datatypes...

...but datatype research is difficult!

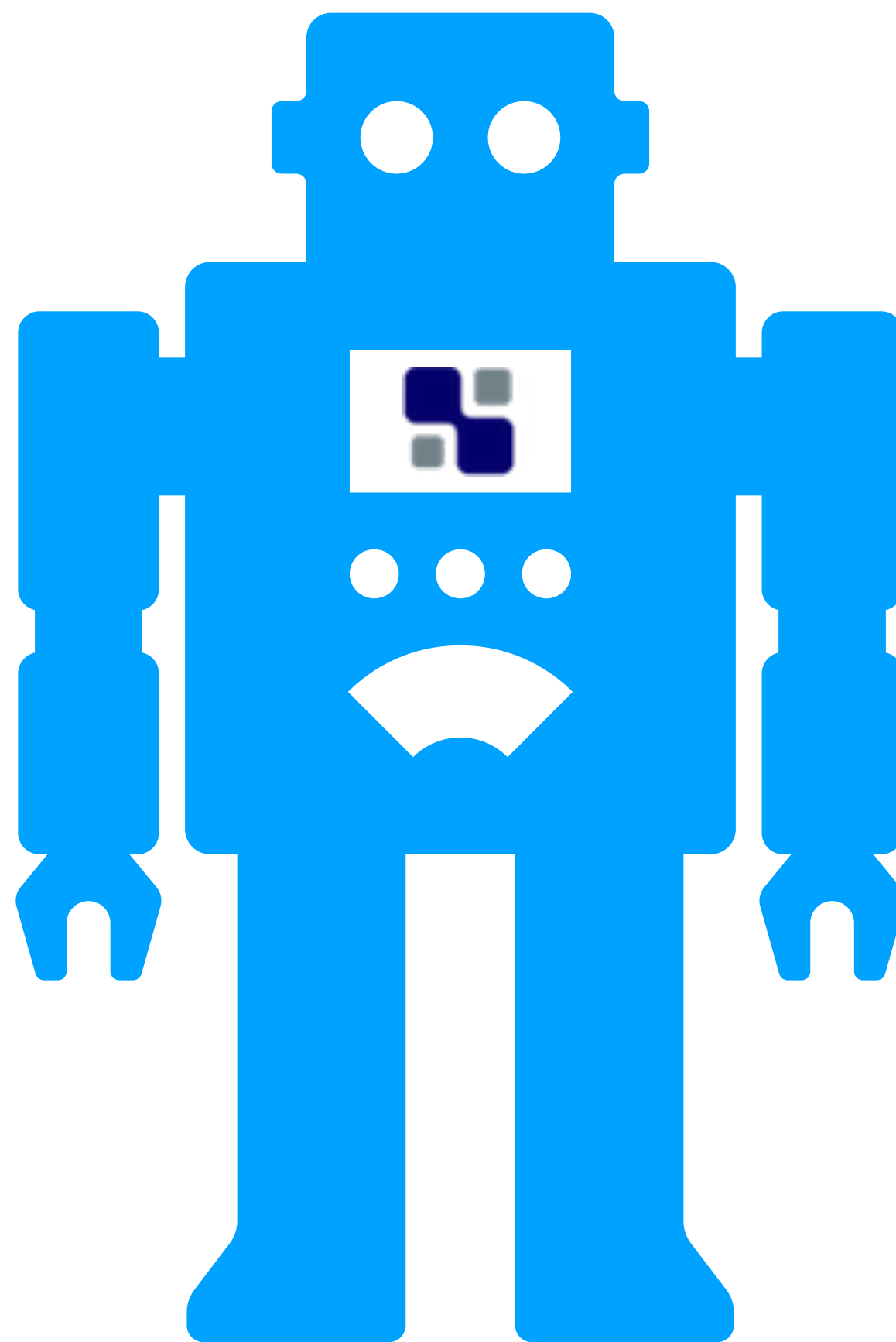
Our solution: the Bring Your Own Datatypes framework.



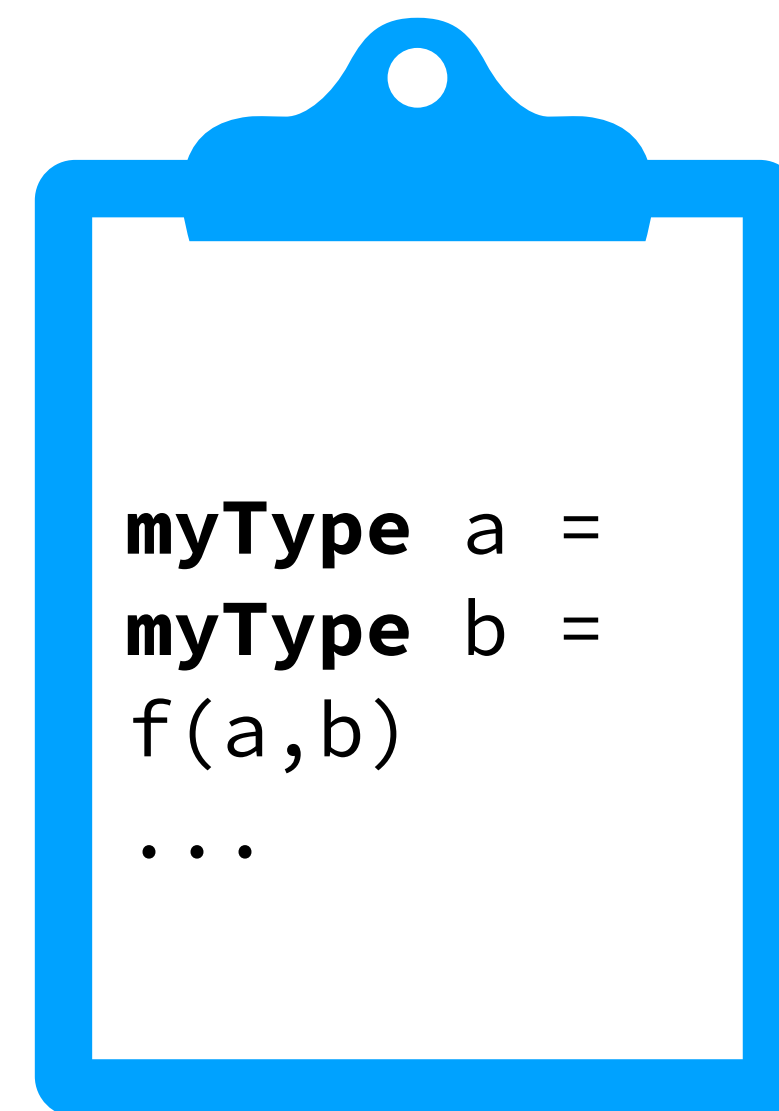
 **tvm**



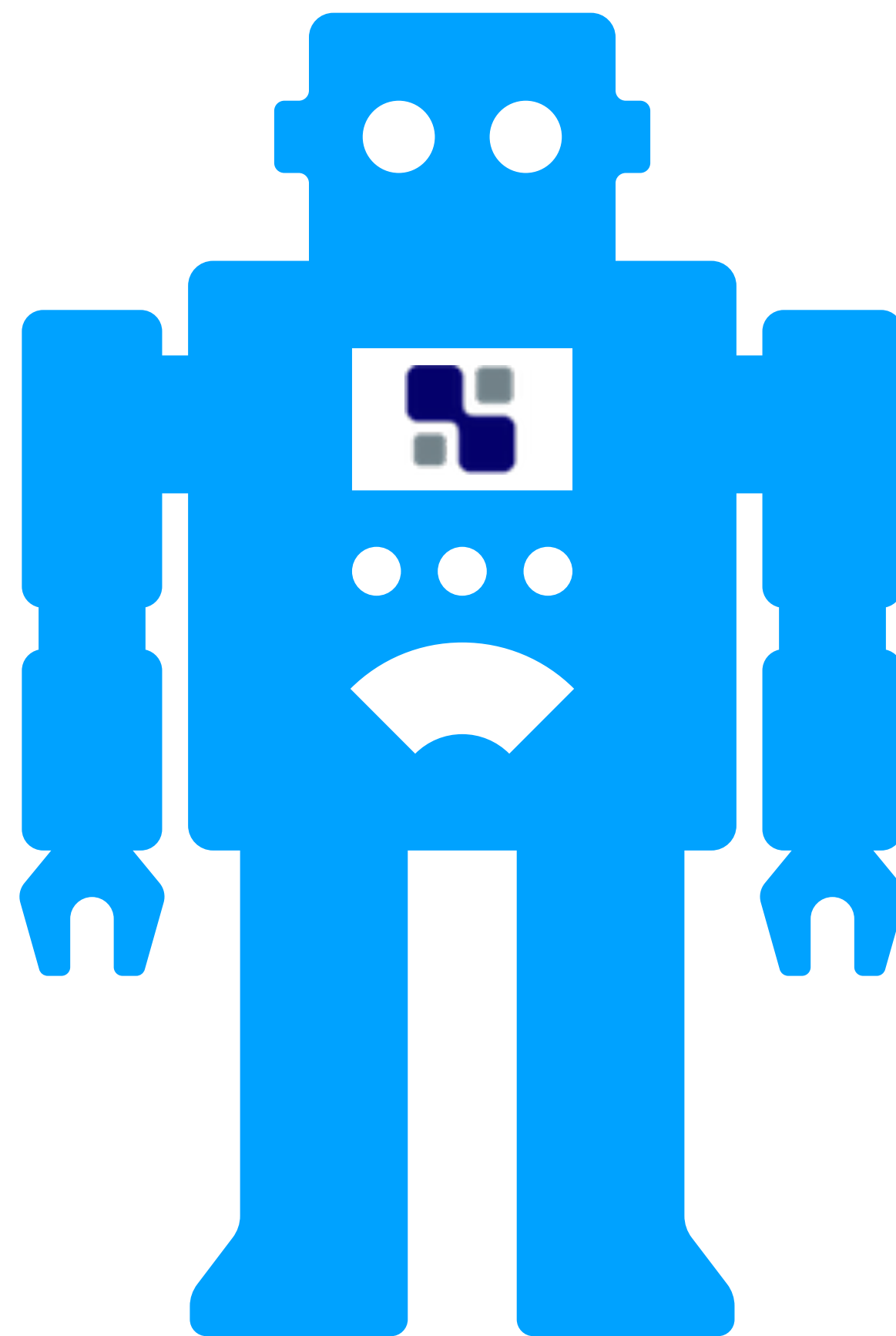
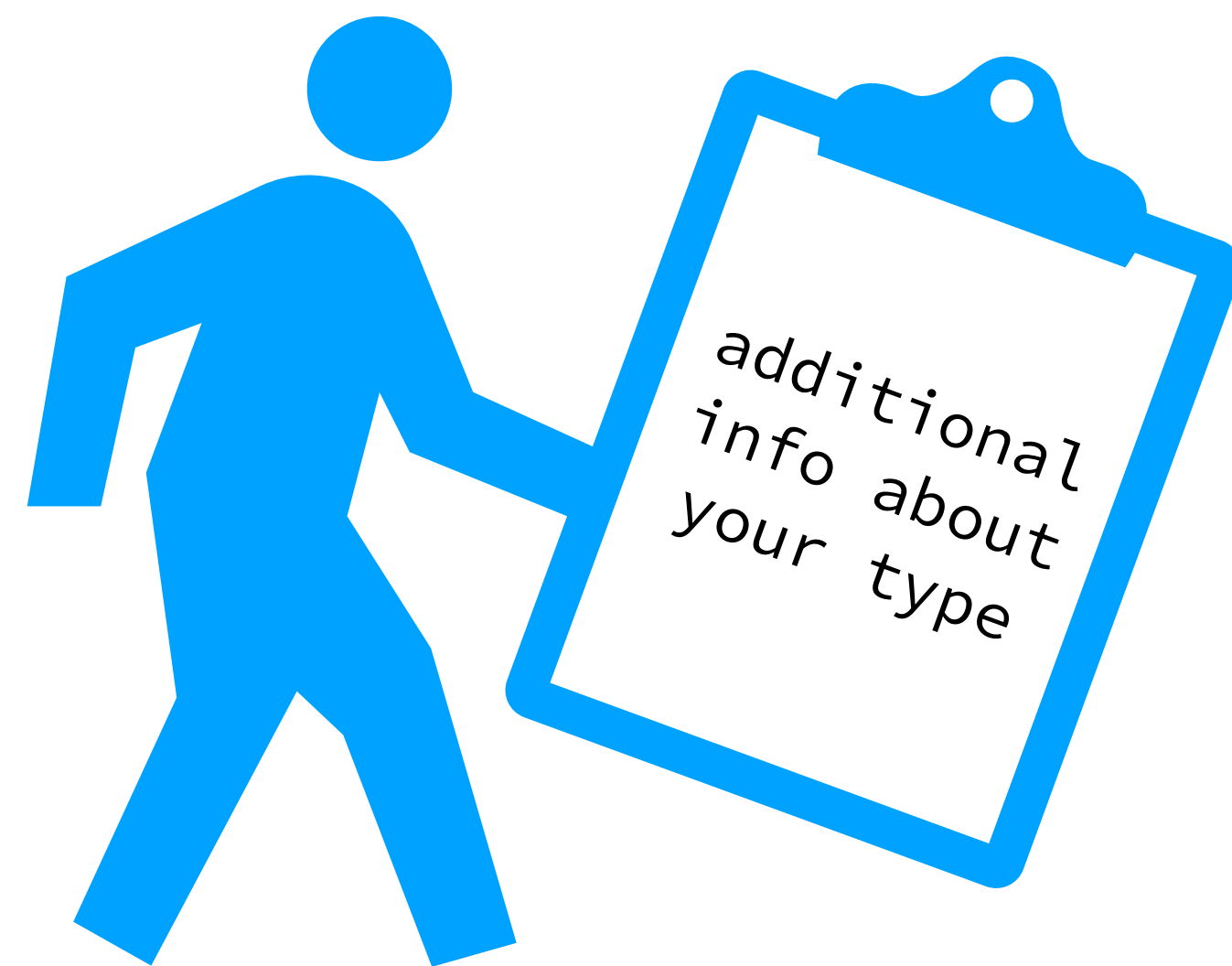
```
myType a =  
myType b =  
f(a,b)  
...
```



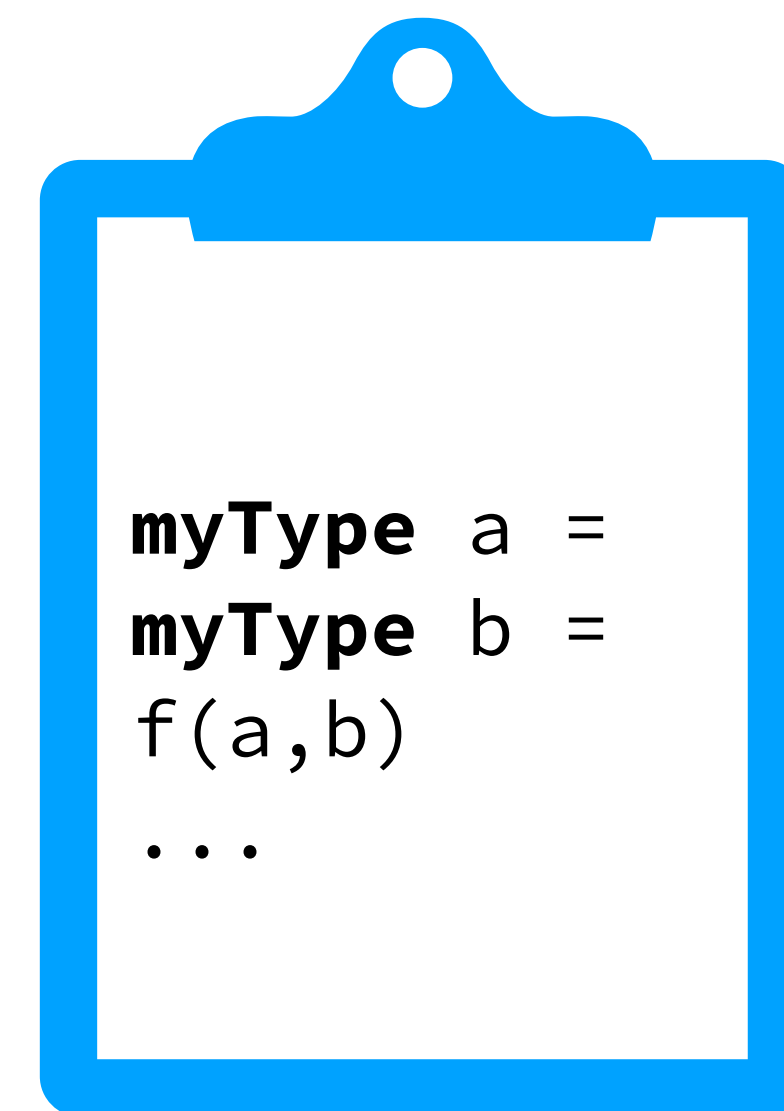
 **tvm**

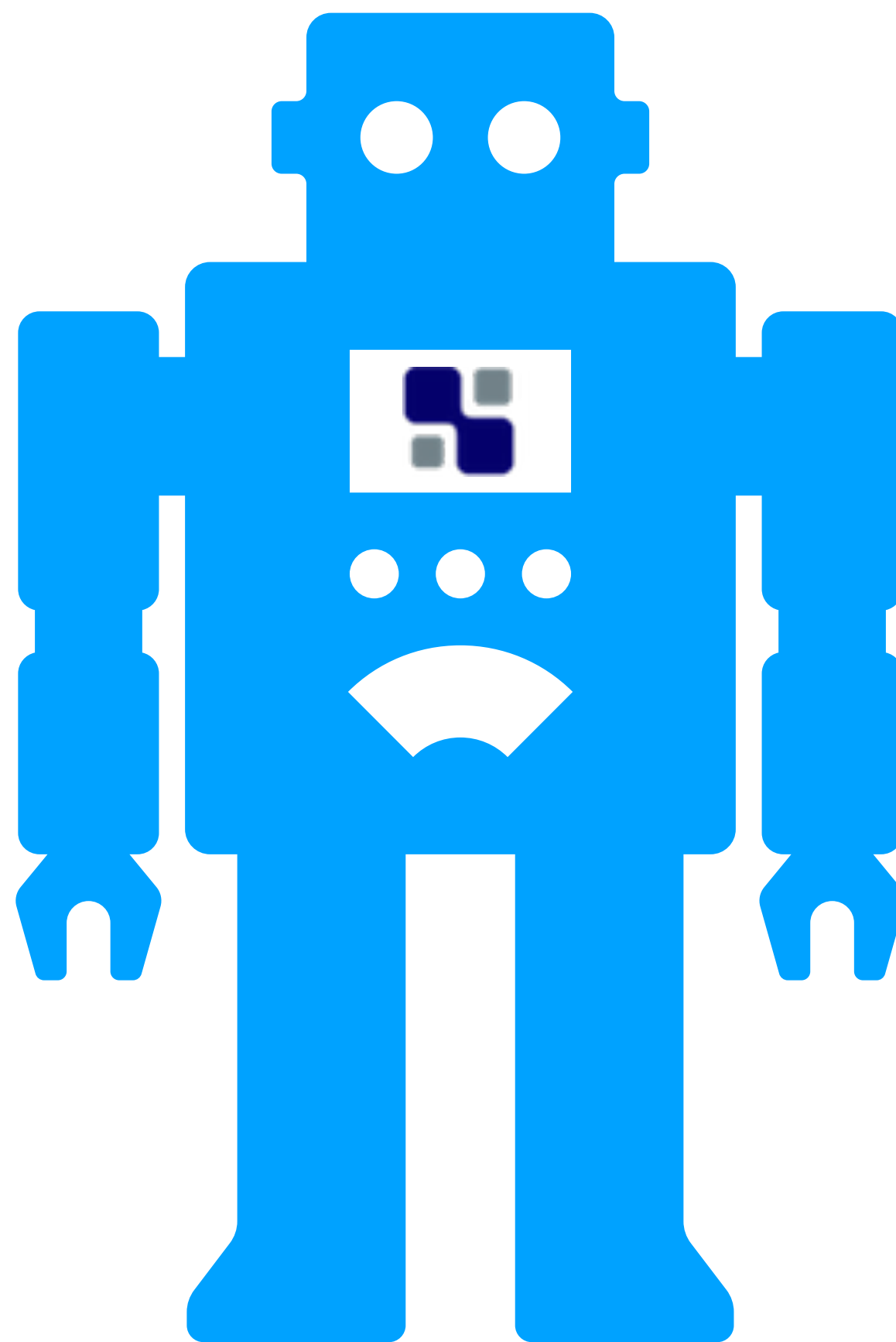


```
myType a =  
myType b =  
f(a,b)  
...
```

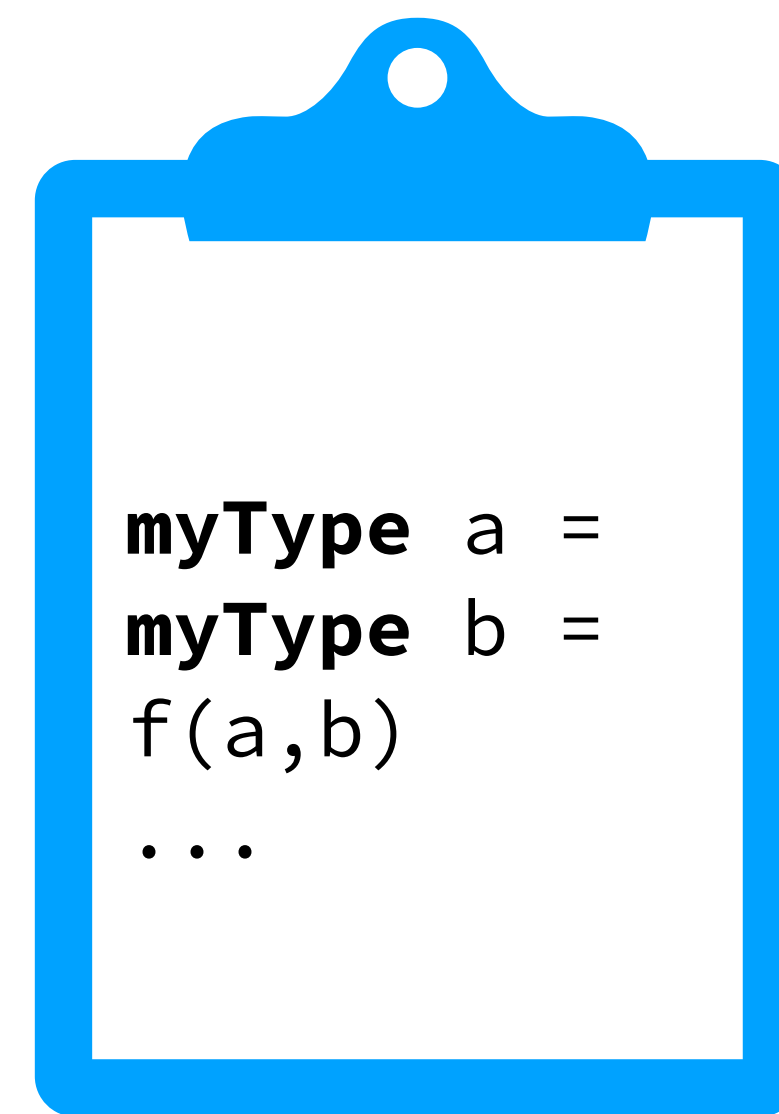



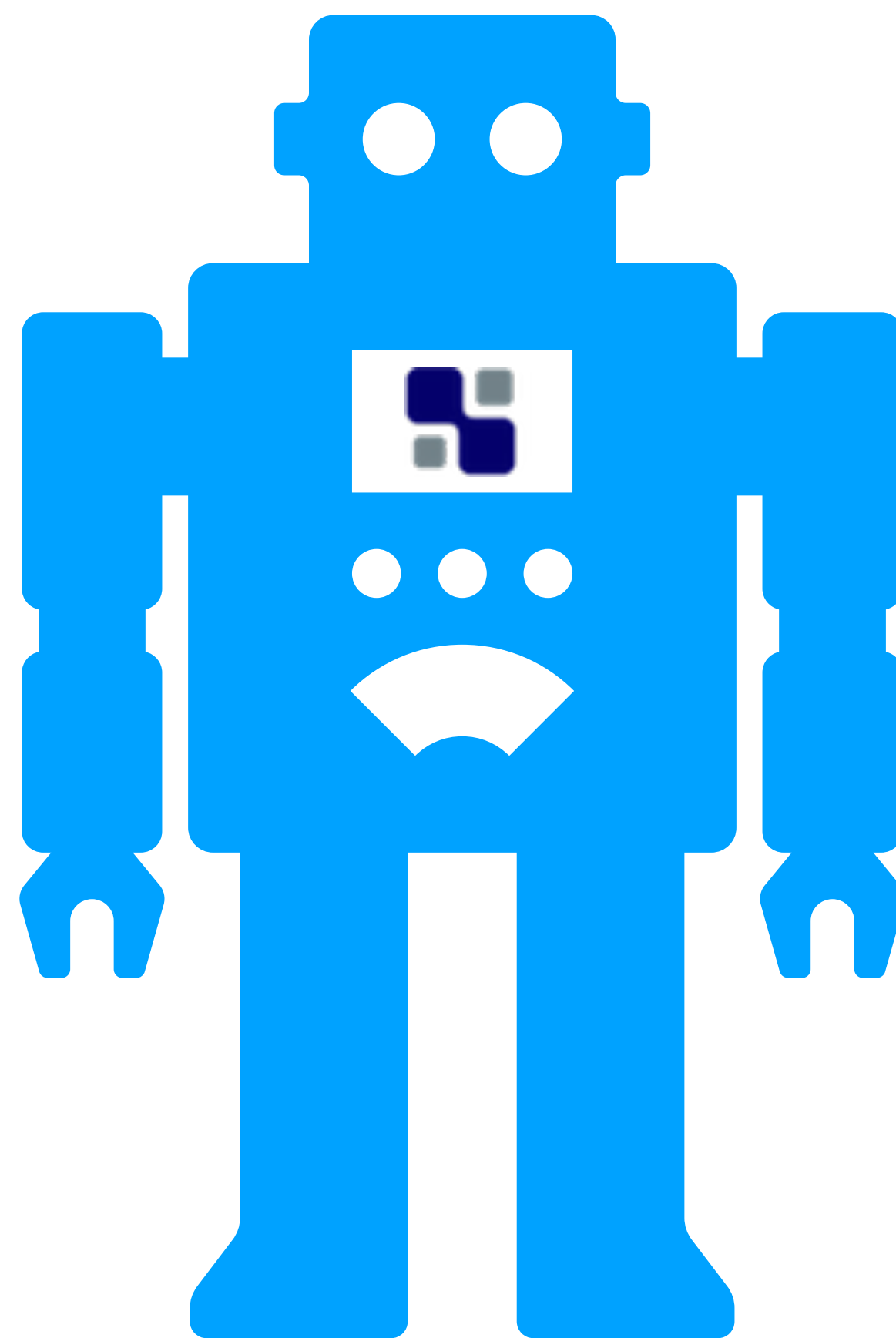
tvm



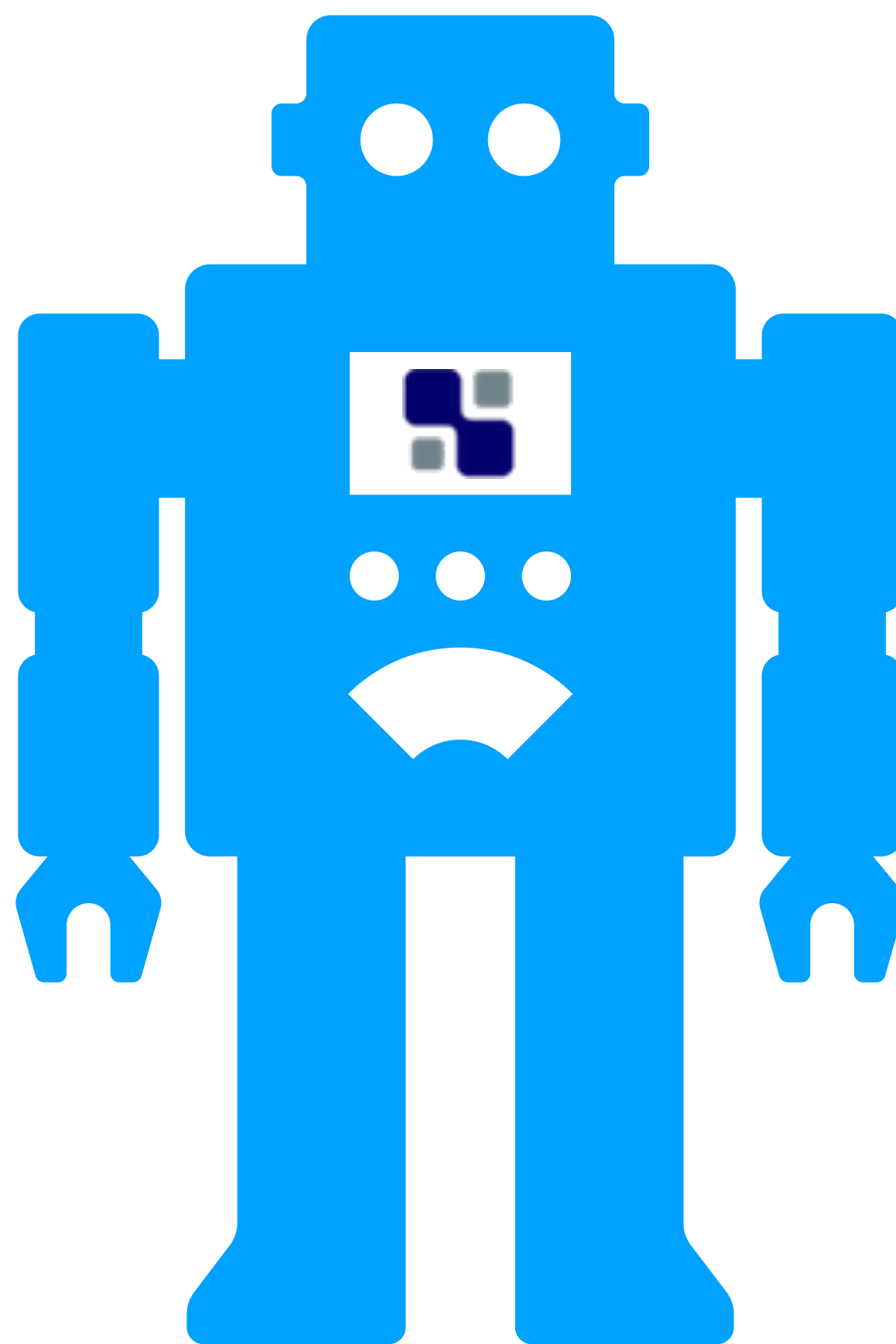


 **tvm**





 tvvm



Bring Your Own Datatypes

What do we want?

What do we want?

- I. User **makes or finds** a custom datatype library which they'd like to use in deep learning workloads

What do we want?

1. User **makes or finds** a custom datatype library which they'd like to use in deep learning workloads
2. User **gives TVM some information** about the library

What do we want?

1. User **makes or finds** a custom datatype library which they'd like to use in deep learning workloads
2. User **gives TVM some information** about the library
3. TVM **compiles and runs programs** which use the custom datatype, handling the custom datatype by calling into the provided library

But first...what is TVM?



But first...what is TVM?



An **extensible, optimizing compiler** for deep learning.

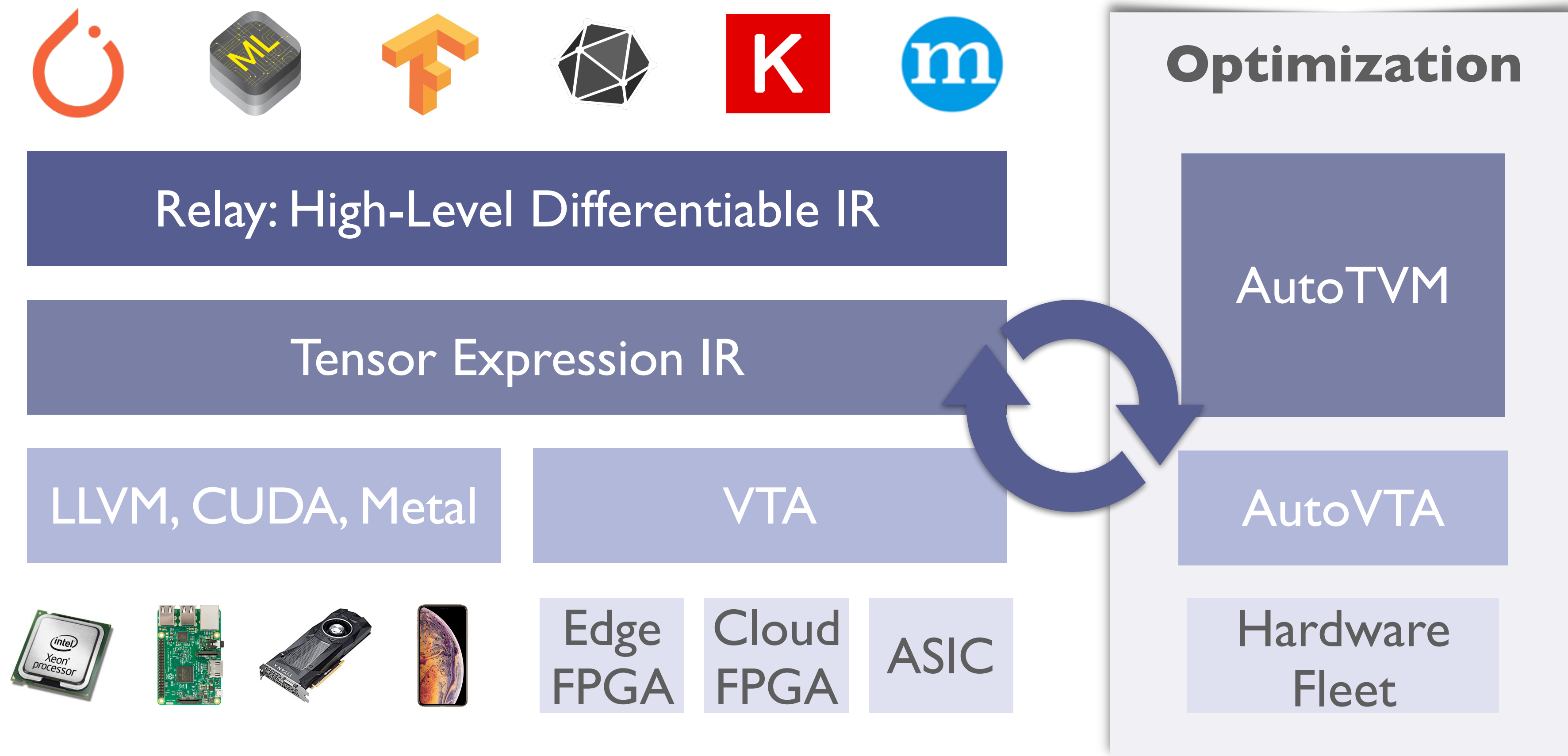
But first...what is TVM?



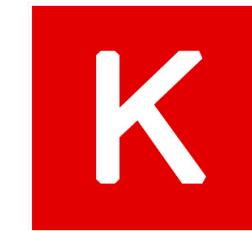
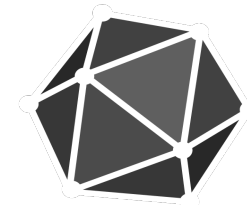
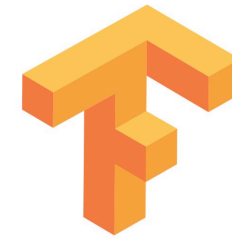
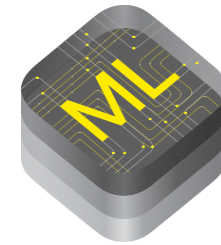
An **extensible, optimizing compiler** for deep learning.

See tvm.ai for more info!

The TVM Stack



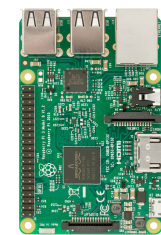
The TVM Stack



TVM IR

LLVM, CUDA, Metal

VTA

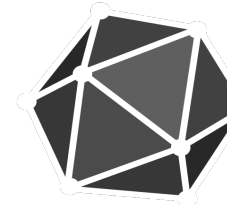
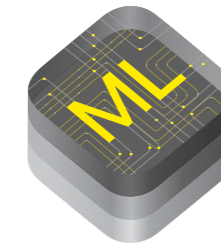


Edge
FPGA

Cloud
FPGA

ASIC

Datatypes Registry



TVM IR

LLVM, CUDA, Metal

VTA

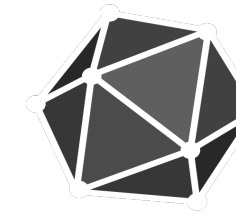
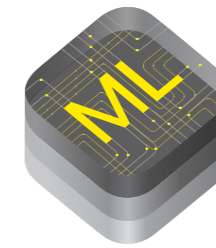


Edge
FPGA

Cloud
FPGA

ASIC

Datatypes Registry



TVM IR

LLVM, CUDA, Metal

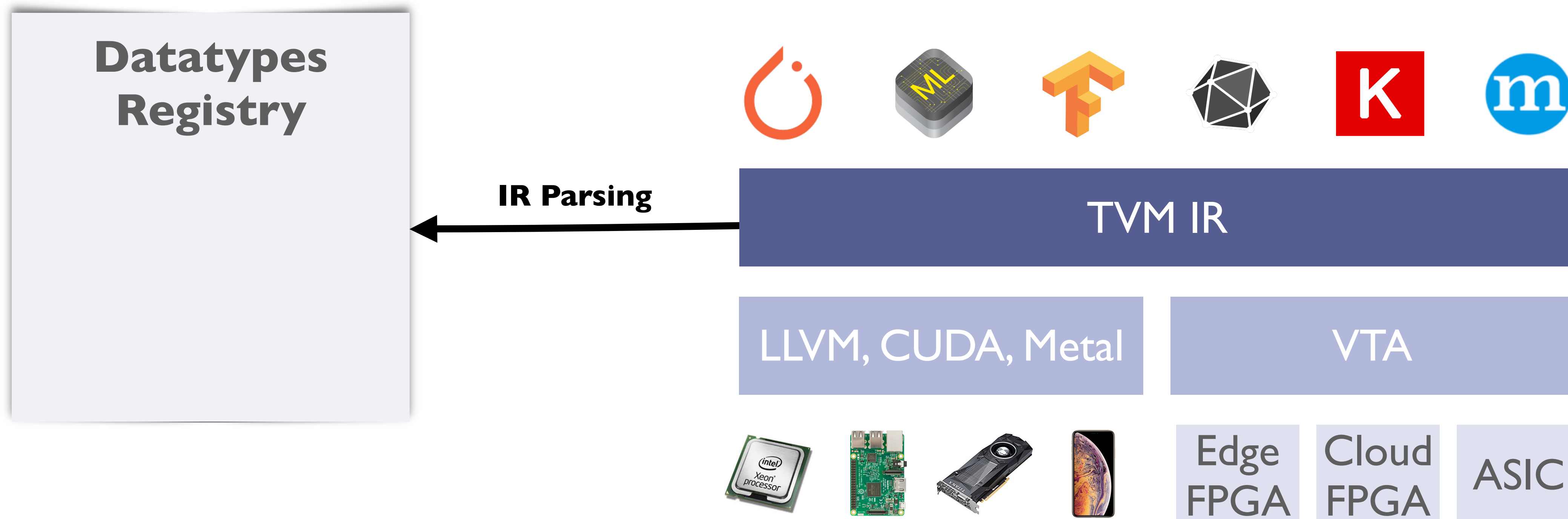
VTA

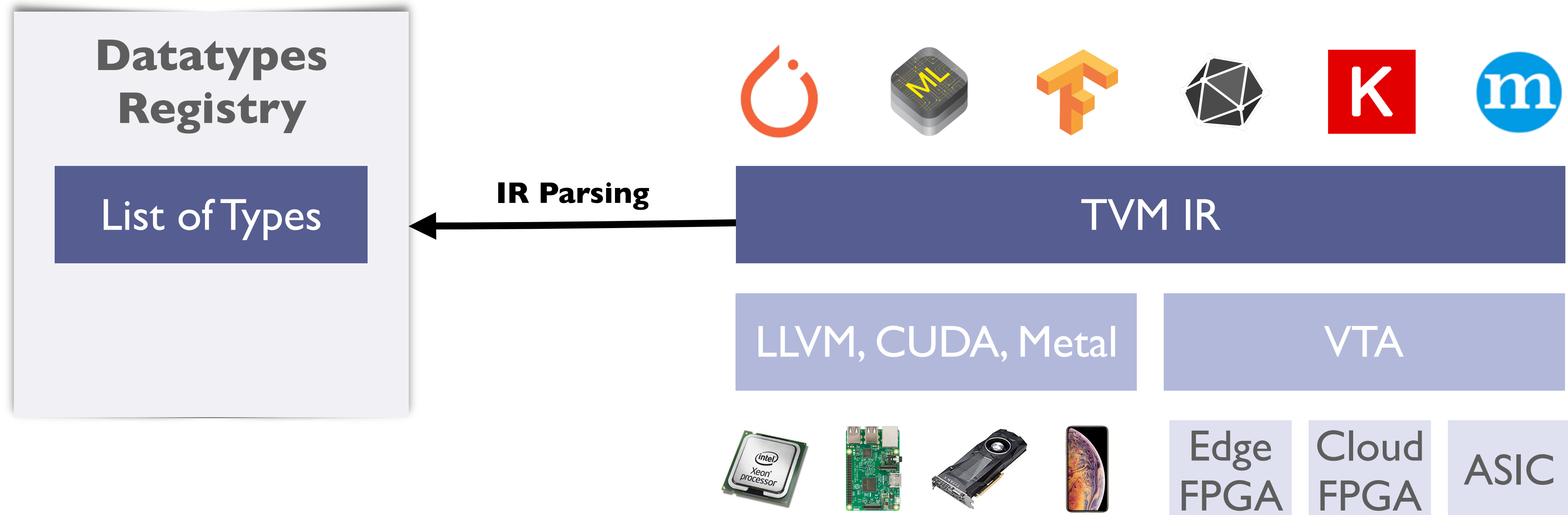


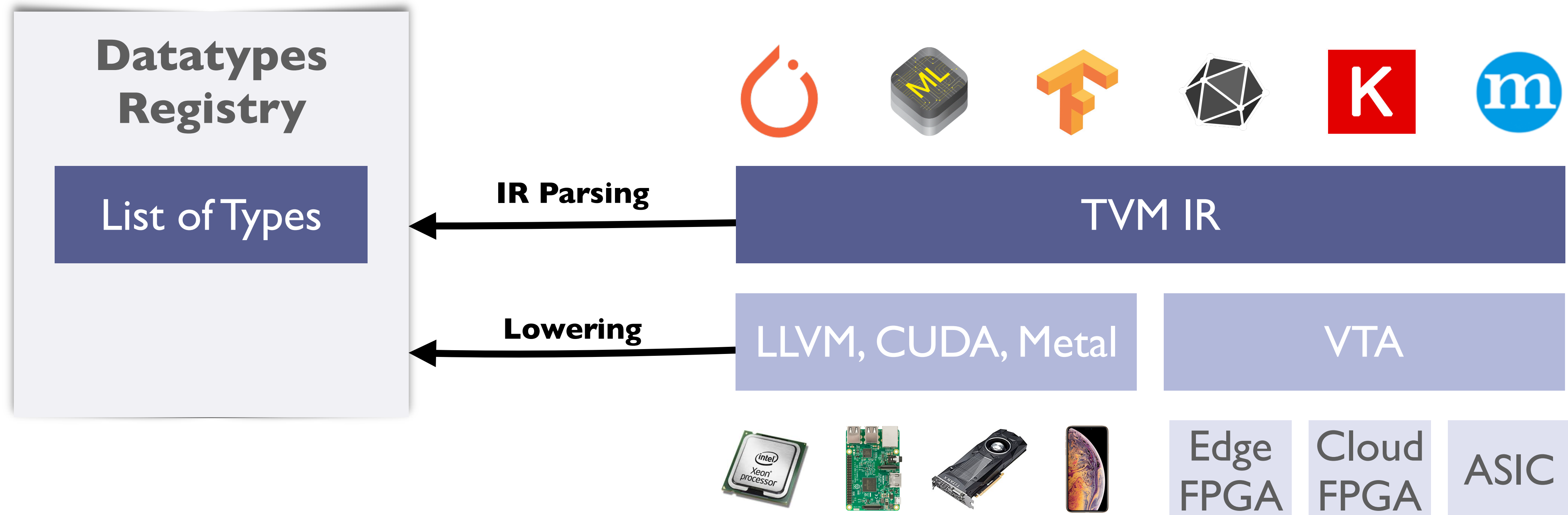
Edge
FPGA

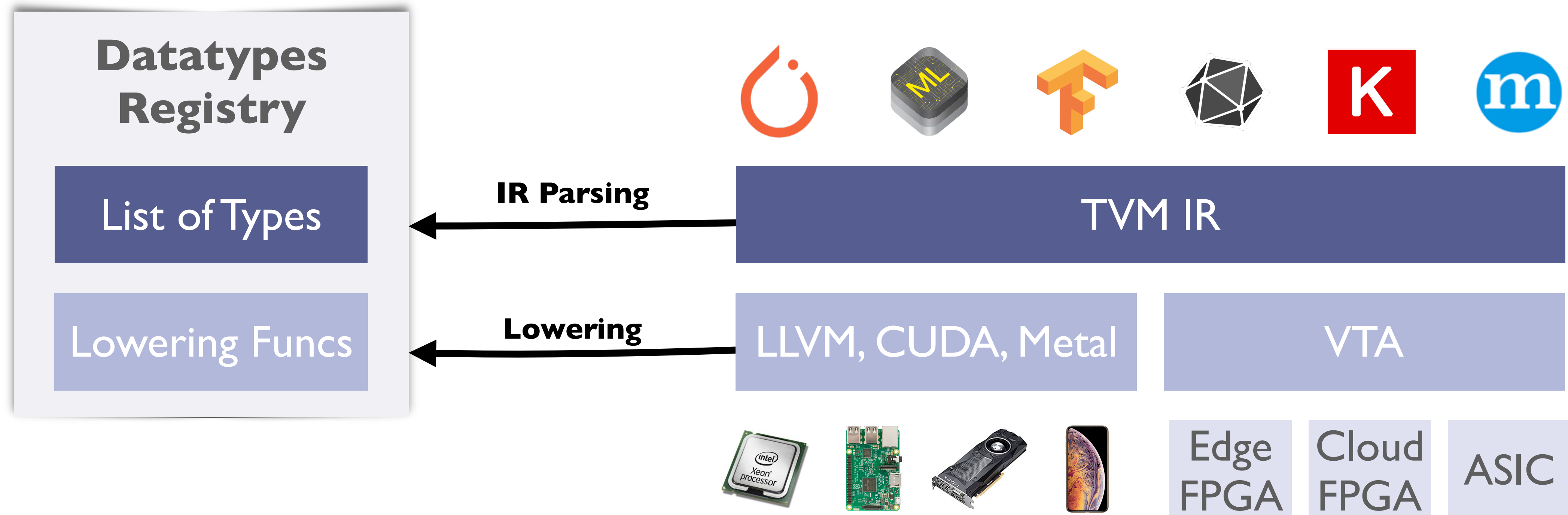
Cloud
FPGA

ASIC









Datatypes Registry

List of Types

Lowering Funcs

Datatypes Registry

List of Types

Lowering Funcs

To register their datatype, the user provides a name and type code:

Datatypes Registry

List of Types

Lowering Funcs

To register their datatype, the user provides a name and type code:

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


Datatypes Registry

List of Types

Lowering Funcs

To register their datatype, the user provides a name and type code:

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

Datatypes Registry

List of Types

Lowering Funcs

To register their datatype, the user provides a name and type code:

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



Datatypes Registry

List of Types

Lowering Funcs

To register their datatype, the user provides a name and type code:

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



This allows TVM to parse programs which use the datatype!

Datatypes Registry

List of Types

Lowering Funcs

Datatypes Registry

List of Types

Lowering Funcs

Users register lowering functions with a similar API:

Datatypes Registry

List of Types

Lowering Funcs

Users register lowering functions with a similar API:

Datatypes Registry

List of Types

Lowering Funcs

Users register lowering functions with a similar API:

```
tvm.datatype.register_op(
```

Datatypes Registry

List of Types

Lowering Funcs

Users register lowering functions with a similar API:

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

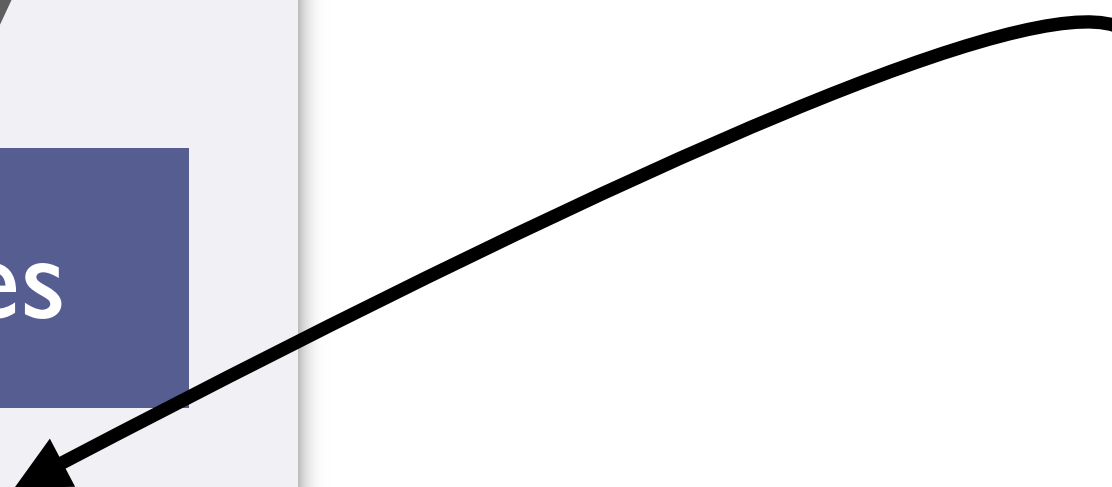

Users register lowering functions with a similar API:

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

Datatypes Registry

List of Types

Lowering Funcs



Datatypes Registry

List of Types

Lowering Funcs

Users register lowering functions with a similar API:

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

This registers a lowering function which lowers bfloat Adds when compiling to LLVM.

Datatypes Registry

List of Types

Lowering Funcs

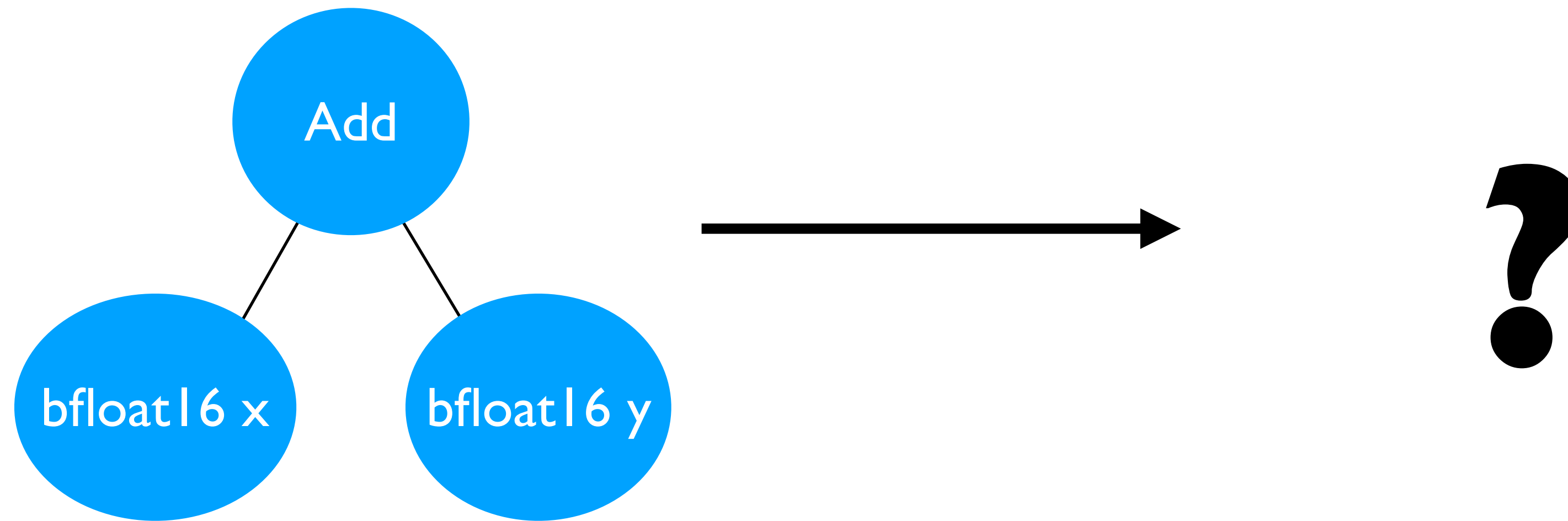
Users register lowering functions with a similar API:

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

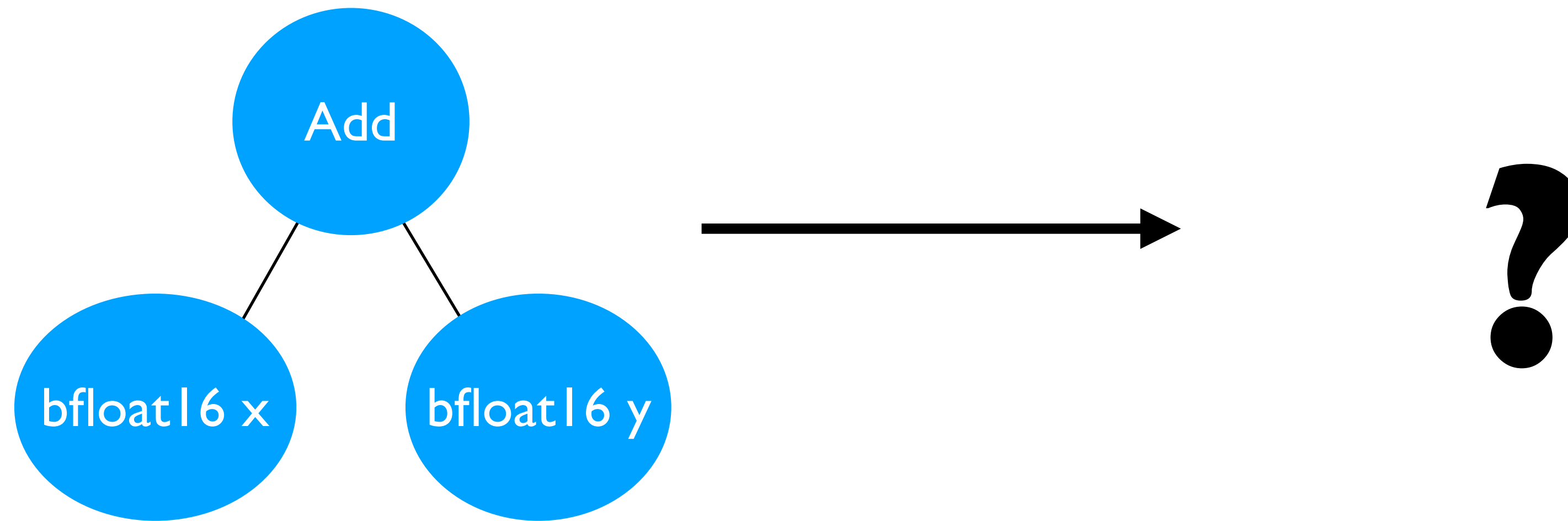
This registers a lowering function which lowers bfloat Adds when compiling to LLVM.

TVM will later use this lowering function when it spits out code!

What do lowering functions look like?

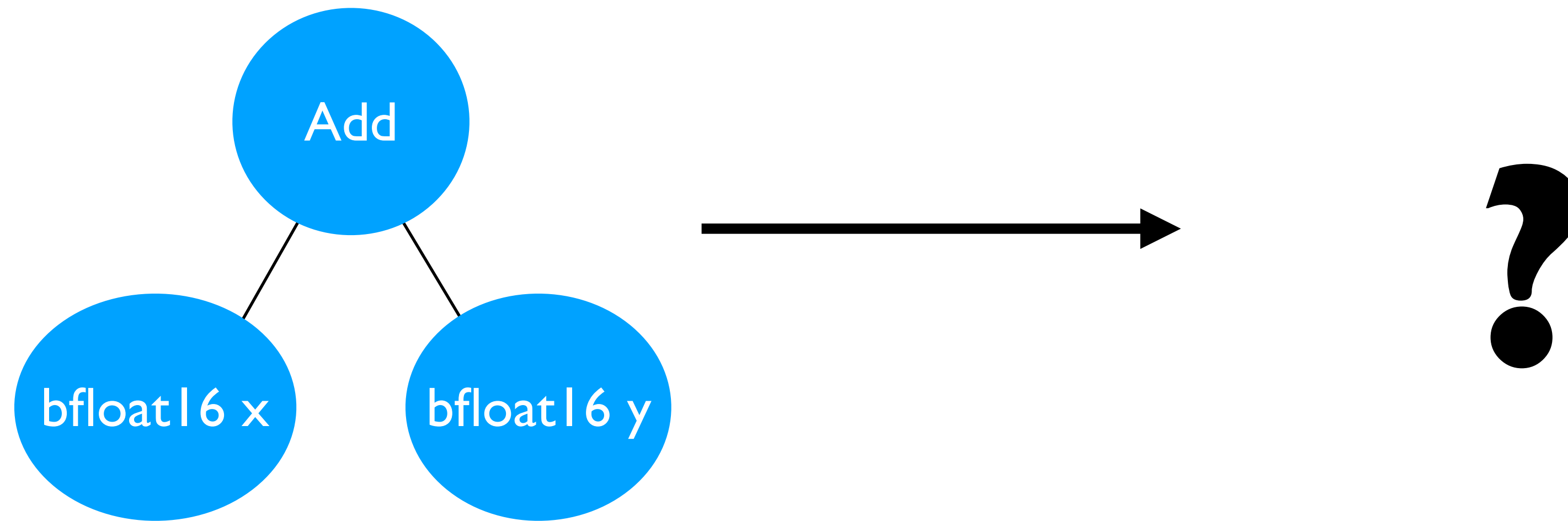


What do lowering functions look like?



Lowering functions convert programs using custom datatypes into programs that native TVM can understand and compile.

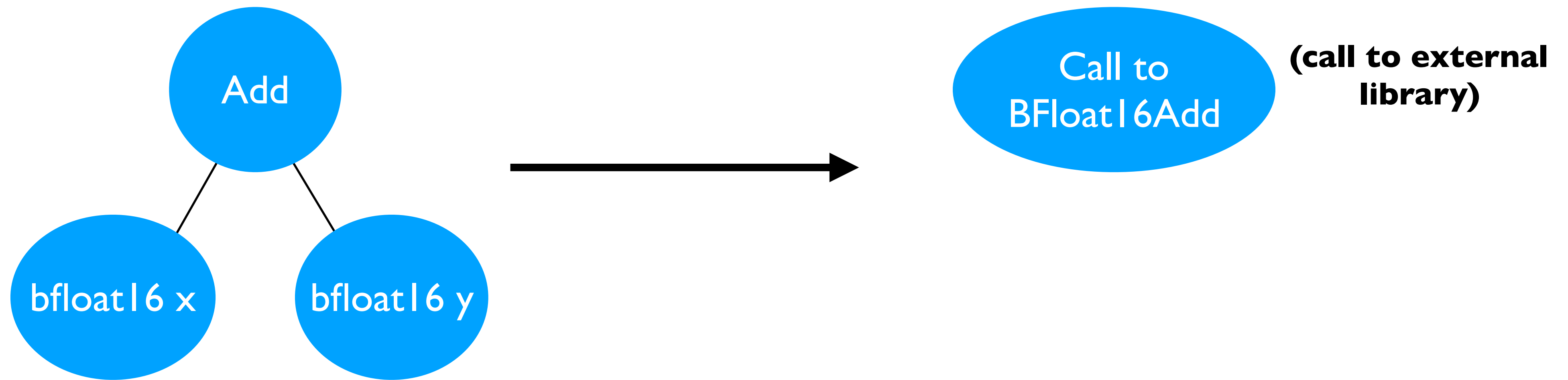
What do lowering functions look like?



Lowering functions convert programs using custom datatypes into programs that native TVM can understand and compile.

In our case, we know our Add over bfloats should become a call to our bfloat library!

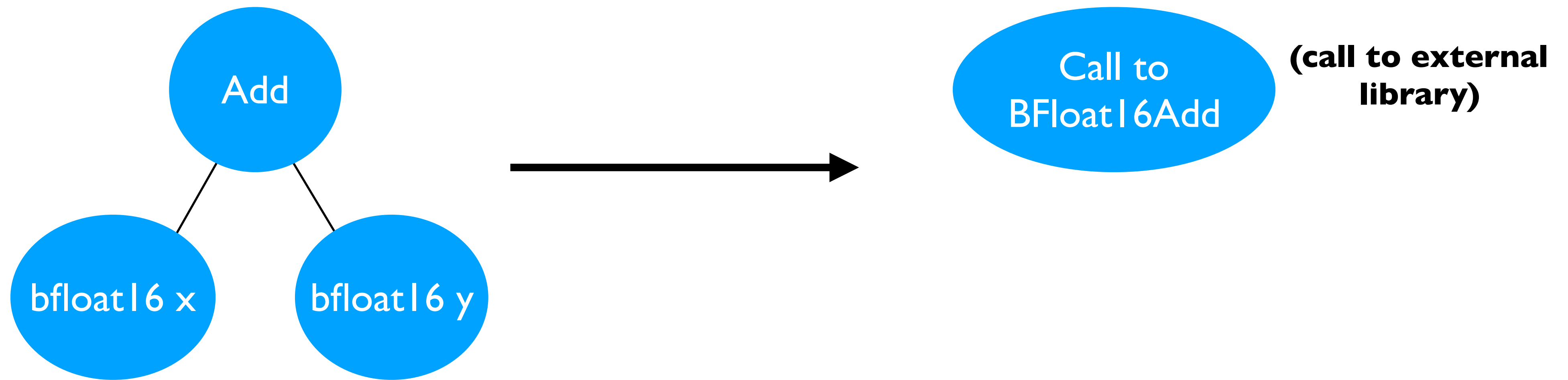
What do lowering functions look like?



Lowering functions convert programs using custom datatypes into programs that native TVM can understand and compile.

In our case, we know our Add over bfloats should become a call to our bfloat library!

What do lowering functions look like?

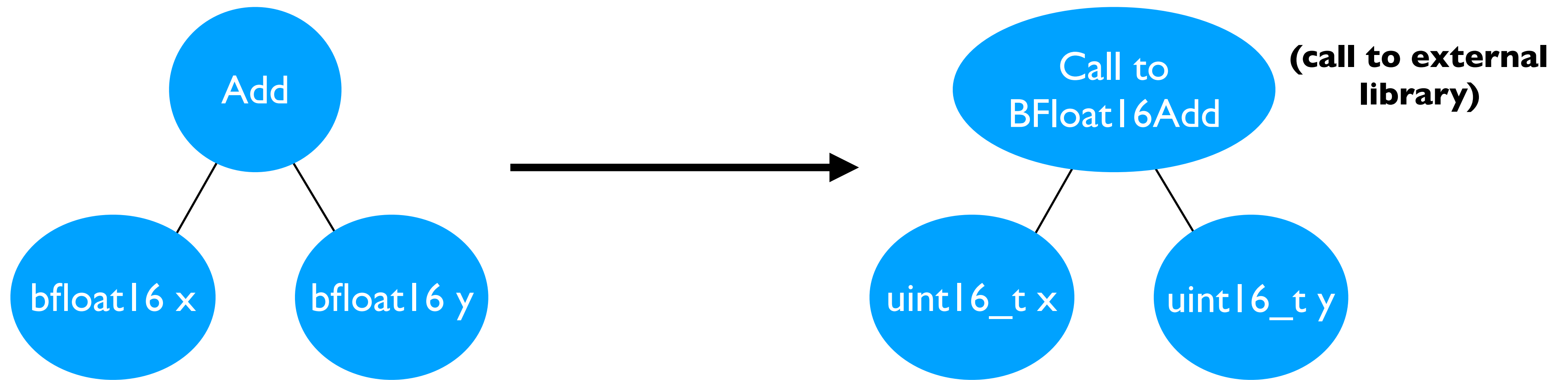


Lowering functions convert programs using custom datatypes into programs that native TVM can understand and compile.

In our case, we know our Add over bfloats should become a call to our bfloat library!

And the inputs? They can just be hidden in opaque unsigned integer types!

What do lowering functions look like?

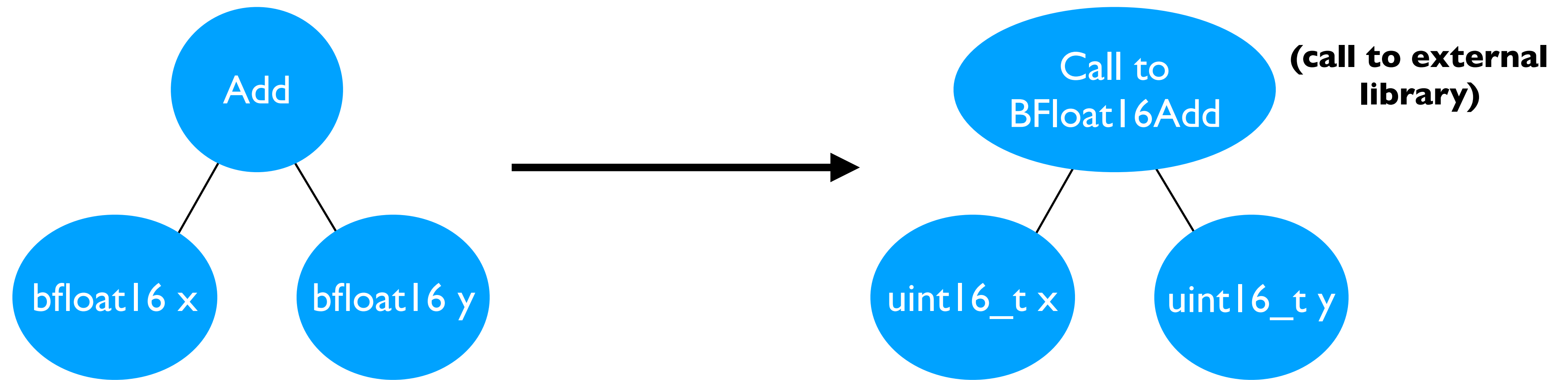


Lowering functions convert programs using custom datatypes into programs that native TVM can understand and compile.

In our case, we know our Add over bfloats should become a call to our bfloat library!

And the inputs? They can just be hidden in opaque unsigned integer types!

What do lowering functions look like?



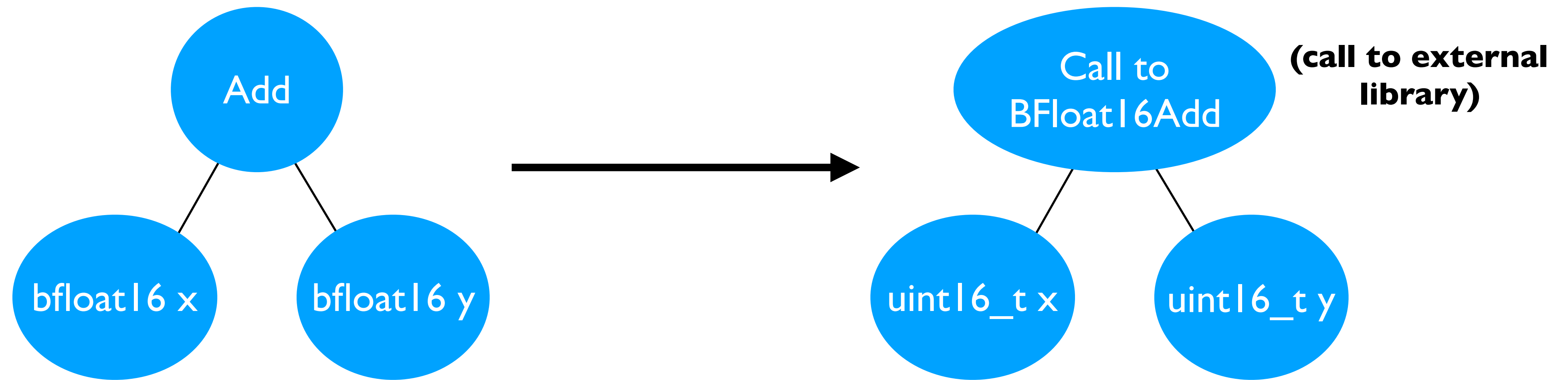
Lowering functions convert programs using custom datatypes into programs that native TVM can understand and compile.

In our case, we know our Add over bfloats should become a call to our bfloat library!

And the inputs? They can just be hidden in opaque unsigned integer types!

We provide a helper function for creating this type of lowering function:

What do lowering functions look like?



Lowering functions convert programs using custom datatypes into programs that native TVM can understand and compile.

In our case, we know our Add over bfloats should become a call to our bfloat library!

And the inputs? They can just be hidden in opaque unsigned integer types!

We provide a helper function for creating this type of lowering function:

```
tvm.datatype.create_lower_func("BFloat16Add")
```

What do we want?

1. User **makes or finds** a custom datatype library which they'd like to use in deep learning workloads
2. User **gives TVM some information** about the library
3. TVM **compiles and runs programs** which use the custom datatype, handling the custom datatype by calling into the provided library

What do we have?

1. User **makes or finds** a custom datatype library which they'd like to use in deep learning workloads
2. User **gives TVM some information** about the library
3. TVM **compiles and runs programs** which use the custom datatype, handling the custom datatype by calling into the provided library

What do we have?

1. User **makes or finds** a custom datatype library which they'd like to use in deep learning workloads
2. User **gives TVM some information** about the library
 - Datatype name
3. TVM **compiles and runs programs** which use the custom datatype, handling the custom datatype by calling into the provided library

What do we have?

1. User **makes or finds** a custom datatype library which they'd like to use in deep learning workloads
2. User **gives TVM some information** about the library
 - Datatype name
 - Lowering functions—user just provides names of library functions!
3. TVM **compiles and runs programs** which use the custom datatype, handling the custom datatype by calling into the provided library

Evaluation

To exercise the framework, I decided to conduct a preliminary evaluation of how a model's trained accuracy changes as we change the datatype.

To exercise the framework, I decided to conduct a preliminary evaluation of how a model's trained accuracy changes as we change the datatype.

We will first discuss the experiment itself and its results—then, we will reflect on the utility of the framework.

Experiment Design

Experiment Design

I. Gathered a list of datatypes

Experiment Design

- I. Gathered a list of datatypes
 - TVM-native

Experiment Design

I. Gathered a list of datatypes

- TVM-native
- Hand-made

Experiment Design

I. Gathered a list of datatypes

- TVM-native
- Hand-made
- From GitHub

Experiment Design

1. Gathered a list of datatypes
 - TVM-native
 - Hand-made
 - From GitHub
2. Pretrained models on the entire CIFAR-10 training set (50k images, 10 classes) in float32 using PyTorch

Experiment Design

1. Gathered a list of datatypes
 - TVM-native
 - Hand-made
 - From GitHub
2. Pretrained models on the entire CIFAR-10 training set (50k images, 10 classes) in float32 using PyTorch
3. Converted the pretrained weights to the custom datatypes.

Experiment Design

1. Gathered a list of datatypes
 - TVM-native
 - Hand-made
 - From GitHub
2. Pretrained models on the entire CIFAR-10 training set (50k images, 10 classes) in float32 using PyTorch
3. Converted the pretrained weights to the custom datatypes.
 - This was done without retraining

Experiment Design

1. Gathered a list of datatypes
 - TVM-native
 - Hand-made
 - From GitHub
2. Pretrained models on the entire CIFAR-10 training set (50k images, 10 classes) in float32 using PyTorch
3. Converted the pretrained weights to the custom datatypes.
 - This was done without retraining
4. Ran the models with the converted datatypes over a sample of the CIFAR-10 test set (100 images each) using TVM

Datatypes:

Datatypes:

- **TVM-native float32** (not using the framework)

Datatypes:

- **TVM-native float32** (not using the framework)
- **float32** (using the framework)

Datatypes:

- **TVM-native float32** (not using the framework)
- **float32** (using the framework)
- Two implementations of **posit8es0**, **posit16es1**, **posit32es2**:
 - Stillwater Supercomputing's Universal library
 - libposit

Datatypes:

- **TVM-native float32** (not using the framework)
- **float32** (using the framework)
- Two implementations of **posit8es0**, **posit16es1**, **posit32es2**:
 - Stillwater Supercomputing's Universal library
 - libposit
- Two implementations of **bfloat16**:
 - My own “naive” implementation
 - biovault-bfloat16: another implementation from GitHub

Datatypes:

- **TVM-native float32** (not using the framework)
- **float32** (using the framework)
- Two implementations of **posit8es0**, **posit16es1**, **posit32es2**:
 - Stillwater Supercomputing's Universal library
 - libposit
- Two implementations of **bfloat16**:
 - My own “naive” implementation
 - biovault-bfloat16: another implementation from GitHub
- “**noptype**”: always returns 0

Models:

- MobilenetV1
- Resnet50

Experiment Results and Evaluation

	resnet accuracy	mobilenet accuracy
float32	0.77	0.71
our bfloat16	0.08	0.11
biovault bfloat16	0.1	0.1
Universal posit8	0.08	0.1
Universal posit16	0.77	0.71
Universal posit32	0.77	0.71
libposit posit8	0.08	0.1
libposit posit16	0.77	0.71
libposit posit32	0.77	0.71

	resnet accuracy	mobilenet accuracy
float32	0.77	0.71
our bfloat16	0.08	0.11
biovault bfloat16	0.1	0.1
Universal posit8	0.08	0.1
Universal posit16	0.77	0.71
Universal posit32	0.77	0.71
libposit posit8	0.08	0.1
libposit posit16	0.77	0.71
libposit posit32	0.77	0.71

	resnet accuracy	mobilenet accuracy
float32	0.77	0.71
our bfloat16	0.08	0.11
biovault bfloat16	0.1	0.1
Universal posit8	0.08	0.1
Universal posit16	0.77	0.71
Universal posit32	0.77	0.71
libposit posit8	0.08	0.1
libposit posit16	0.77	0.71
libposit posit32	0.77	0.71

	resnet accuracy	mobilenet accuracy
float32	0.77	0.71
our bfloat16	0.08	0.11
biovault bfloat16	0.1	0.1
Universal posit8	0.08	0.1
Universal posit16	0.77	0.71
Universal posit32	0.77	0.71
libposit posit8	0.08	0.1
libposit posit16	0.77	0.71
libposit posit32	0.77	0.71

	resnet accuracy	mobilenet accuracy
float32	0.77	0.71
our bfloat16	0.08	0.11
biovault bfloat16	0.1	0.1
Universal posit8	0.08	0.1
Universal posit16	0.77	0.71
Universal posit32	0.77	0.71
libposit posit8	0.08	0.1
libposit posit16	0.77	0.71
libposit posit32	0.77	0.71

	resnet accuracy	mobilenet accuracy
float32	0.77	0.71
our bfloat16	0.08	0.11
biovault bfloat16	0.1	0.1
Universal posit8	0.08	0.1
✓ Universal posit16	0.77	0.71
✓ Universal posit32	0.77	0.71
libposit posit8	0.08	0.1
✓ libposit posit16	0.77	0.71
✓ libposit posit32	0.77	0.71

	resnet accuracy	mobilenet accuracy
float32	0.77	0.71
our bfloat16	0.08	0.11
biovault bfloat16	0.1	0.1
✗ Universal posit8	0.08	0.1
Universal posit16	0.77	0.71
Universal posit32	0.77	0.71
✗ libposit posit8	0.08	0.1
libposit posit16	0.77	0.71
libposit posit32	0.77	0.71

		resnet accuracy	mobilenet accuracy
	float32	0.77	0.71
Same size!	✗ our bfloat16	0.08	0.11
	✗ biovault bfloat16	0.1	0.1
	Universal posit8	0.08	0.1
	✓ Universal posit16	0.77	0.71
	Universal posit32	0.77	0.71
	libposit posit8	0.08	0.1
	✓ libposit posit16	0.77	0.71
	libposit posit32	0.77	0.71

Framework Evaluation

I wanted to evaluate three aspects of the framework:

I wanted to evaluate three aspects of the framework:

- Overhead

I wanted to evaluate three aspects of the framework:

- Overhead
- Ease of use

I wanted to evaluate three aspects of the framework:

- Overhead
- Ease of use
- Breadth of datatypes which can be used

Overhead

To measure overhead, I compared the inference time of three types:

To measure overhead, I compared the inference time of three types:

- TVM-native float32

To measure overhead, I compared the inference time of three types:

- TVM-native float32
- float32 implemented in the framework

To measure overhead, I compared the inference time of three types:

- TVM-native float32
- float32 implemented in the framework
- “noptype”: a custom type that does no work

mobilenet v1

	mean inference time (s)	framework overhead
native float32	0.10	1x
float32	0.12	1.11x
noptype	0.11	1.10x

resnet50

	mean inference time (s)	framework overhead
native float32	0.21	1x
float32	0.52	2.51x
noptype	0.28	1.35x

mobilenet v1

	mean inference time (s)	framework overhead
native float32	0.10	1x
float32	0.12	1.11x
noptype	0.11	1.10x

resnet50

	mean inference time (s)	framework overhead
native float32	0.21	1x
float32	0.52	2.51x
noptype	0.28	1.35x

mobilenet v1

	mean inference time (s)	framework overhead
native float32	0.10	1x
float32	0.12	1.11x
noptype	0.11	1.10x

resnet50

	mean inference time (s)	framework overhead
native float32	0.21	1x
float32	0.52	2.51x
noptype	0.28	1.35x

mobilenet v1

	mean inference time (s)	framework overhead
native float32	0.10	1x
float32	0.12	1.11x
noptype	0.11	1.10x

resnet50

	mean inference time (s)	framework overhead
native float32	0.21	1x
float32	0.52	2.51x
noptype	0.28	1.35x

mobilenet v1

	mean inference time (s)	framework overhead
native float32	0.10	1x
float32	0.12	1.11x
noptype	0.11	1.10x

resnet50

	mean inference time (s)	framework overhead
native float32	0.21	1x
float32	0.52	2.51x
noptype	0.28	1.35x

Main source of overhead: not in added computation, but in the **compilation opportunity cost**.

mobilenet v1

mean inference time (s) framework overhead

native float32

0.10

1x

float32

0.12

1.11x

noptype

0.11

1.10x

resnet50

mean inference time (s) framework overhead

native float32

0.21

1x

float32

0.52

2.51x

noptype

0.28

1.35x

Main source of overhead: not in added computation, but in the **compilation opportunity cost**.

- noptype is fairly low-overhead in both workloads → function calls don't add too much

mobilenet v1

	mean inference time (s)	framework overhead
--	-------------------------	--------------------

native float32	0.10	1x
----------------	------	----

resnet50

	mean inference time (s)	framework overhead
--	-------------------------	--------------------

native float32	0.21	1x
----------------	------	----

float32	0.12	1.11x
---------	------	-------

float32	0.52	2.51x
---------	------	-------

noptype	0.11	1.10x
---------	------	-------

noptype	0.28	1.35x
---------	------	-------

Main source of overhead: not in added computation, but in the **compilation opportunity cost**.

- noptype is fairly low-overhead in both workloads → function calls don't add too much
- float32 is low-overhead in Mobilenet, which is optimized for low total float ops

mobilenet v1

	mean inference time (s)	framework overhead
--	-------------------------	--------------------

native float32	0.10	1x
----------------	------	----

float32	0.12	1.11x
---------	------	-------

noptype	0.11	1.10x
---------	------	-------

resnet50

	mean inference time (s)	framework overhead
--	-------------------------	--------------------

native float32	0.21	1x
----------------	------	----

float32	0.52	2.51x
---------	------	-------

noptype	0.28	1.35x
---------	------	-------

Main source of overhead: not in added computation, but in the **compilation opportunity cost**.

- noptype is fairly low-overhead in both workloads → function calls don't add too much
- float32 is low-overhead in Mobilenet, which is optimized for low total float ops
- ...but float32 is high-overhead in ResNet → native float32 ResNet much more optimized!

Ease of Use

By the Numbers...

By the Numbers...

To use libposit posit32 in Mobilenet and ResNet...

By the Numbers...

To use libposit posit32 in Mobilenet and ResNet...

- A total of **12 operators** needed to be implemented, taking **about 70 lines of C++ in a wrapper library**.

By the Numbers...

To use libposit posit32 in Mobilenet and ResNet...

- A total of **12 operators** needed to be implemented, taking **about 70 lines of C++ in a wrapper library**.
- Operators including **casts to/from posits, add/sub/mul/div**, more complex math operators like **exp**, and comparators like **max**.

By the Numbers...

To use libposit posit32 in Mobilenet and ResNet...

- A total of **12 operators** needed to be implemented, taking **about 70 lines of C++ in a wrapper library**.
 - Operators including **casts to/from posits, add/sub/mul/div**, more complex math operators like **exp**, and comparators like **max**.
- In the **150-line** Python script which runs ResNet50 with posit32,

By the Numbers...

To use libposit posit32 in Mobilenet and ResNet...

- A total of **12 operators** needed to be implemented, taking **about 70 lines of C++ in a wrapper library**.
 - Operators including **casts to/from posits, add/sub/mul/div**, more complex math operators like **exp**, and comparators like **max**.
- In the **150-line** Python script which runs ResNet50 with posit32,
 - **57 lines** register the datatype and define lowering functions for the 12 operators,

By the Numbers...

To use libposit posit32 in Mobilenet and ResNet...

- A total of **12 operators** needed to be implemented, taking **about 70 lines of C++ in a wrapper library**.
 - Operators including **casts to/from posits, add/sub/mul/div**, more complex math operators like **exp**, and comparators like **max**.
- In the **150-line** Python script which runs ResNet50 with posit32,
 - **57 lines** register the datatype and define lowering functions for the 12 operators,
 - **3 lines** convert the model to posit32,

By the Numbers...

To use libposit posit32 in Mobilenet and ResNet...

- A total of **12 operators** needed to be implemented, taking **about 70 lines of C++ in a wrapper library**.
 - Operators including **casts to/from posits, add/sub/mul/div**, more complex math operators like **exp**, and comparators like **max**.
- In the **150-line** Python script which runs ResNet50 with posit32,
 - **57 lines** register the datatype and define lowering functions for the 12 operators,
 - **3 lines** convert the model to posit32,
 - **3 lines** convert the input, run the model, and convert the output.

How could we improve?

How could we improve?

- Allow the user to specify their own calling convention, removing the need for a wrapper over the library

How could we improve?

- Allow the user to specify their own calling convention, removing the need for a wrapper over the library
- Implement cleaner registration functions in the TVM Python frontend

Breadth of Datatypes

Breadth of Datatypes

Breadth of Datatypes

We were able to successfully represent a small sample of modern datatypes!

Breadth of Datatypes

We were able to successfully represent a small sample of modern datatypes!

Questionable whether current BYOD can represent...

Breadth of Datatypes

We were able to successfully represent a small sample of modern datatypes!

Questionable whether current BYOD can represent...

- Block floating point, or any other type with “external” state

Breadth of Datatypes

We were able to successfully represent a small sample of modern datatypes!

Questionable whether current BYOD can represent...

- Block floating point, or any other type with “external” state
- Datatypes with elements larger than 64 bits

Breadth of Datatypes

We were able to successfully represent a small sample of modern datatypes!

Questionable whether current BYOD can represent...

- Block floating point, or any other type with “external” state
- Datatypes with elements larger than 64 bits

Could potentially be implemented by allowing datatypes to attach metadata to each scalar

Future Work

Future Work

- Training in TVM—ramifications for custom datatypes?

Future Work

- Training in TVM—ramifications for custom datatypes?
- Improve performance

Future Work

- Training in TVM—ramifications for custom datatypes?
- Improve performance
 - Enable inlining of LLVM

Future Work

- Training in TVM—ramifications for custom datatypes?
- Improve performance
 - Enable inlining of LLVM
 - Allow users to supply optimized kernels for higher-level operators, e.g. posit conv2d

Future Work

- Training in TVM—ramifications for custom datatypes?
- Improve performance
 - Enable inlining of LLVM
 - Allow users to supply optimized kernels for higher-level operators, e.g. posit conv2d
- Tackle complex datatypes like block floating point

In conclusion, I have presented the Bring Your Own Datatypes framework.

In conclusion, I have presented the Bring Your Own Datatypes framework.

In conclusion, I have presented the Bring Your Own Datatypes framework.

I have shown how the framework can enable useful datatype research.

Thank You!

Extra Slides

Registering the Datatype

```
1 import tvm
2 from tvm import relay
3 from ctypes import CDLL, RTLD_GLOBAL
4
5 CDLL('libposit-wrapper.so', RTLD_GLOBAL)
6
7 tvm.datatype.register('posit32', 131)
8
9 tvm.datatype.register_op(
10     tvm.datatype.create_lower_func('Posit32Add'),
11     'Add', 'llvm', 'posit32')
```


Registering the Datatype

```
1 import tvm
2 from tvm import relay
3 from ctypes import CDLL, RTLD_GLOBAL
4
5 CDLL('libposit-wrapper.so', RTLD_GLOBAL)
6
7 tvm.datatype.register('posit32', 131)
8
9 tvm.datatype.register_op(
10     tvm.datatype.create_lower_func('Posit32Add'),
11     'Add', 'llvm', 'posit32')
```

✓ Any loaded C-linkage functions will work!



Registering the Datatype

```
1 import tvm
2 from tvm import relay
3 from ctypes import CDLL, RTLD_GLOBAL
4
5 CDLL('libposit-wrapper.so', RTLD_GLOBAL)
6
7 tvm.datatype.register('posit32', 131)
8
9 tvm.datatype.register_op(
10     tvm.datatype.create_lower_func('Posit32Add'),
11     'Add', 'llvm', 'posit32')
```

✓ Any loaded C-linkage functions will work!



✓ Helper functions for common-case

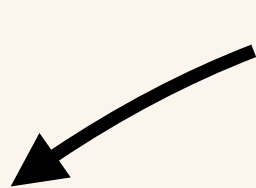


Implementing the Datatype

```
1 #include <stdint>
2 #include <posit.h>
3
4 posit32_t Uint32ToPosit32(uint32_t in) {
5     return posit32_reinterpret(in);
6 }
7
8 uint32_t Posit32ToUint32(posit32_t in) {
9     return posit32_bits(in);
10 }
11
12 extern "C" uint32_t Posit32Add(uint32_t a, uint32_t b) {
13     auto a_p = Uint32ToPosit32(a);
14     auto b_p = Uint32ToPosit32(b);
15     auto sum_p = posit32_add(a_p, b_p);
16     return Posit32ToUint32(sum_p);
17 }
```


Implementing the Datatype

```
1 #include <stdint>
2 #include <posit.h>
3
4 posit32_t Uint32ToPosit32(uint32_t in) {
5     return posit32_reinterpret(in);
6 }
7
8 uint32_t Posit32ToUint32(posit32_t in) {
9     return posit32_bits(in);
10 }
11
12 extern "C" uint32_t Posit32Add(uint32_t a, uint32_t b) {
13     auto a_p = Uint32ToPosit32(a);
14     auto b_p = Uint32ToPosit32(b);
15     auto sum_p = posit32_add(a_p, b_p);
16     return Posit32ToUint32(sum_p);
17 }
```

 **✗ Inflexible calling convention**


Converting the Program

```
1 module, params = load_resnet50()  
2 module, params = change_dtype('float32', 'custom[posit32]32',  
3                               module['main'], params)  
4  
5 ex = relay.create_executor(mod=module)  
6 model = ex.evaluate()
```


Converting the Program

```
1 module, params = load_resnet50()  
2 module, params = change_dtype('float32', 'custom[posit32]32',  
3                               module['main'], params)  
4  
5 ex = relay.create_executor(mod=module)  
6 model = ex.evaluate()
```

✓ Changing datatype of a program is easy!




Running the Program

```
1 image = ... # load image
2
3 image = convert_ndarray('custom[posit32]32', image)
4
5 with tvm.build_config(disable_vectorize=True):
6     output = model(image, **params)
7
8 output = convert_ndarray('float32', output)
```

Running the Program

```
1 image = ... # load image
2
3 image = convert_ndarray('custom[posit32]32', image)
4
5 with tvm.build_config(disable_vectorize=True):
6     output = model(image, **params)
7
8 output = convert_ndarray('float32', output)
```

✓ Converting data is also easy!



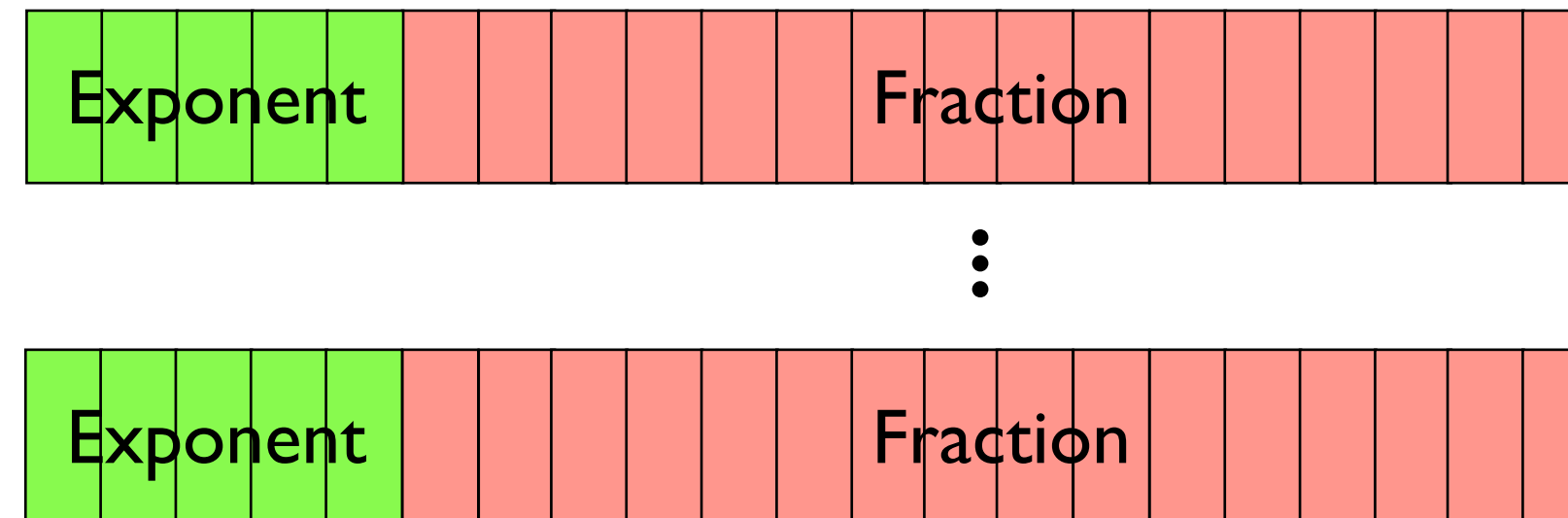
Running the Program

```
1 image = ... # load image
2
3 image = convert_ndarray('custom[posit32]32', image)
4
5 with tvn.build_config(disable_vectorize=True):
6     output = model(image, **params)
7
8 output = convert_ndarray('float32', output)
```

✓ Converting data is also easy!

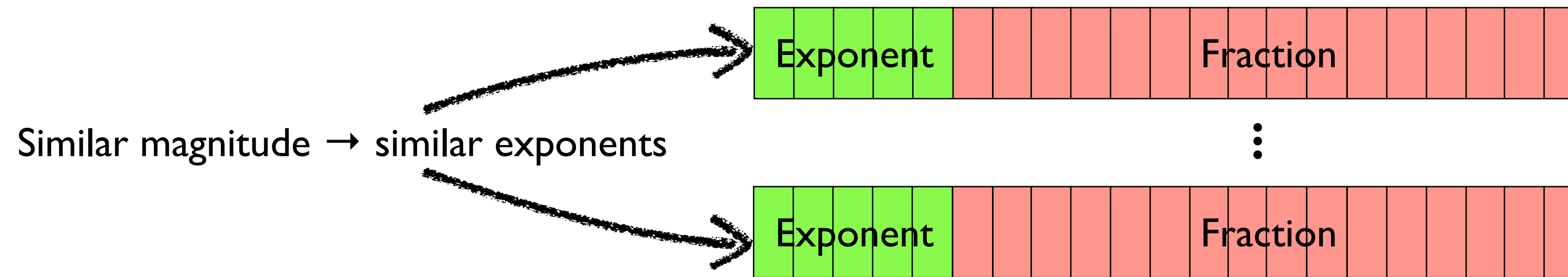
✗ Have to manually disable vectorization

Flexpoint

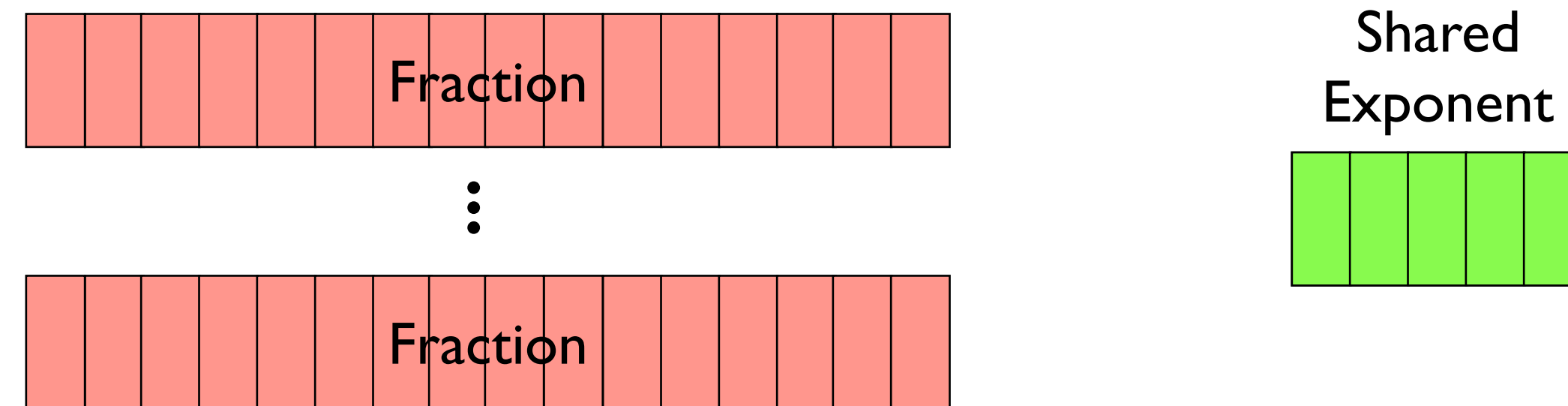


See Köster et. al., [“An Adaptive Numerical Format for Efficient Training of Deep Neural Networks”](#)

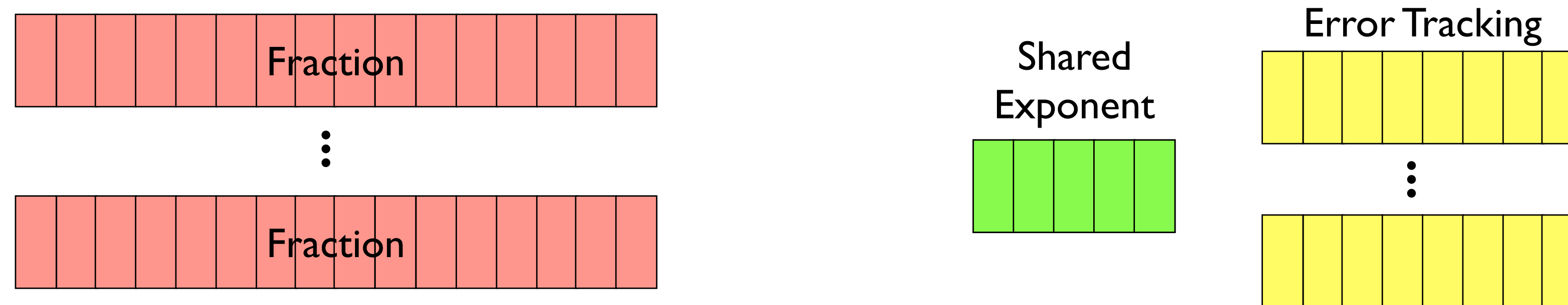
Flexpoint



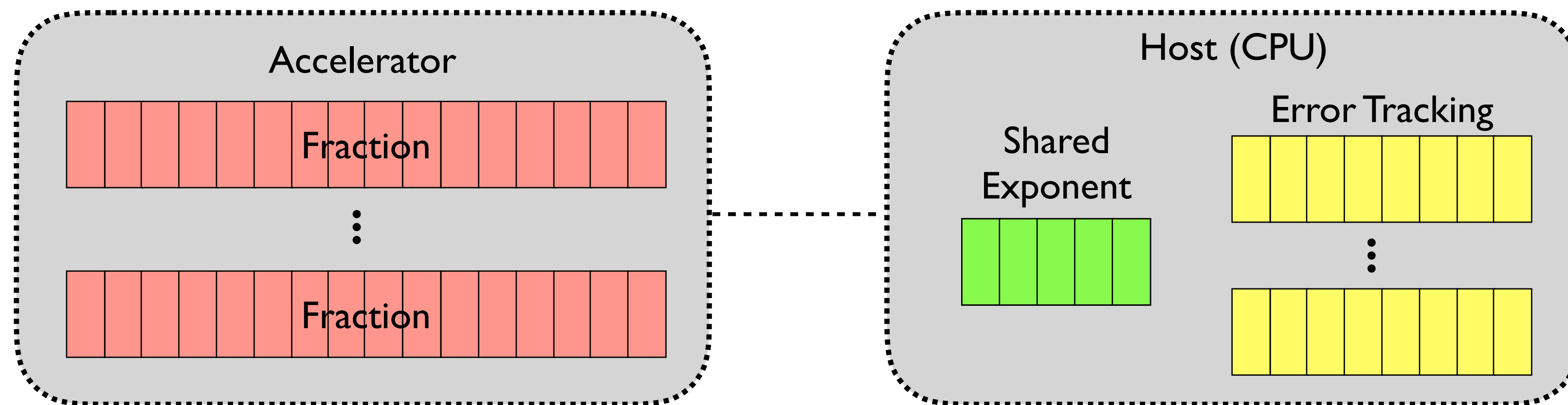
Flexpoint



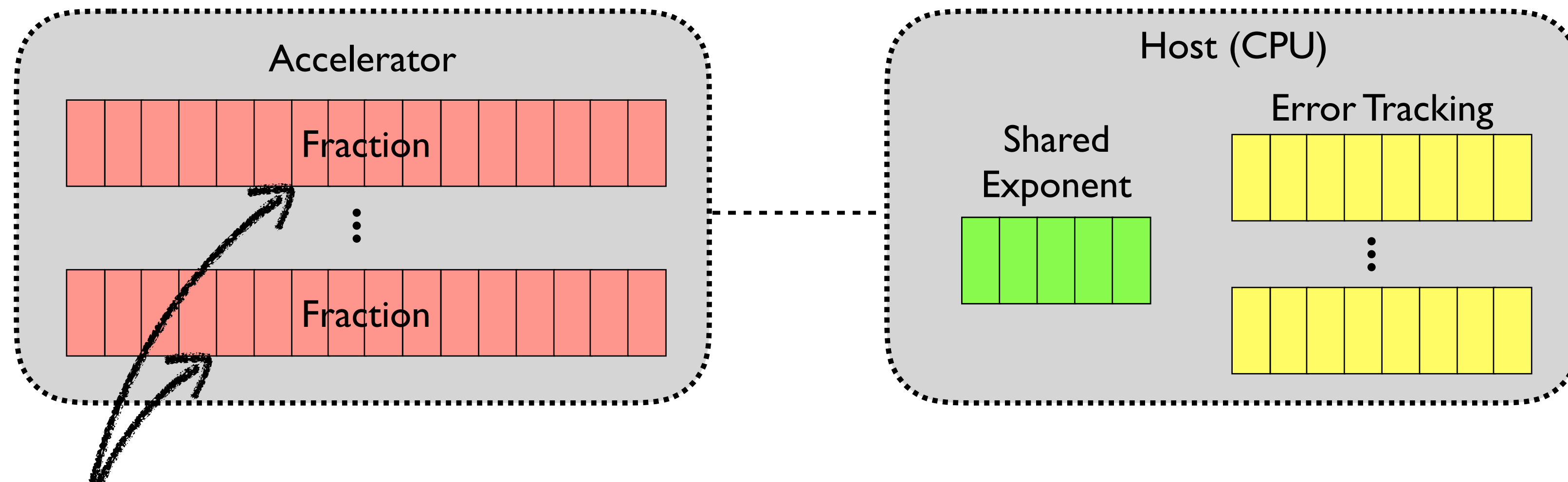
Flexpoint



Flexpoint



Flexpoint



Just integers! We can use integer hardware!