

Juice - Machine Learning with Rust

A brief introduction and architecture overview

Bernhard Schuster

November 30, 2017

Scope of Juice: Deep Learning, Artificial Neural Networks

Scope of Juice: Deep Learning

Does:


- Deep Learning / Neural Networks

Not:

- Classic Machine Learning
- Bayes
- k-Means
- Decision Trees
- Support Vector Machine

What I am not going to explain/show

- network architecture design
- speed comparisons
- fancy graphics

- Leaf
- Autumn AI 
- Goal: Re-pricing prediction for online merchants



1 Juice ← Leaf



2 Greenglas ← Cuticula

- Coaster ← Collenchyma



3 Org: Sparrow

¹make

²me




³pretty

Motivation



Competition: Tensorflow, Torch, Caffe, Theano

Sparrow Goals

- Easy to run
- Easy to integrate
- Easy to deploy on   
- Ready for Embedded ⁴

⁴anything 32bits and up with enough memory

An aerial photograph of a water park. The image shows several slides of different colors: yellow, green, and blue. A red umbrella is open on a paved area in the center. A person is visible on a blue slide on the right side. The park is surrounded by lush green trees and foliage. A dark grey banner with white text is overlaid across the middle of the image.

rusty-machine, vikos, rustlearn, prophet, nn, ...

<your application>

juice

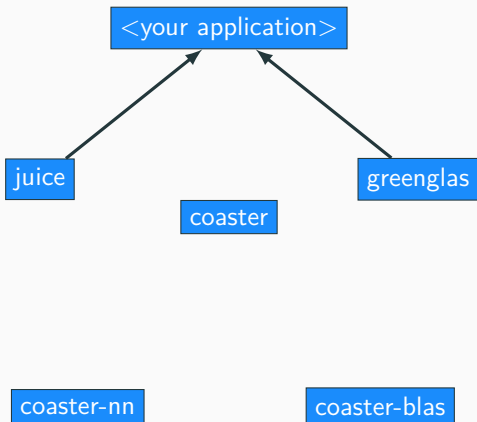
greenglas

coaster

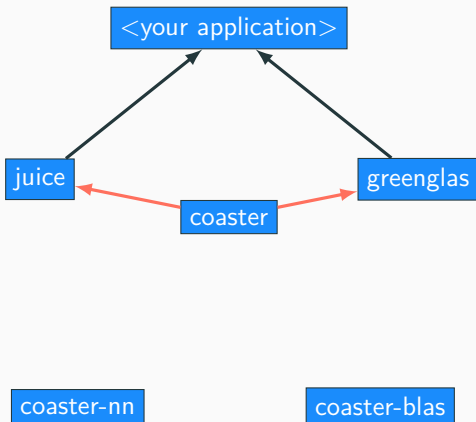
coaster-nn

coaster-blas

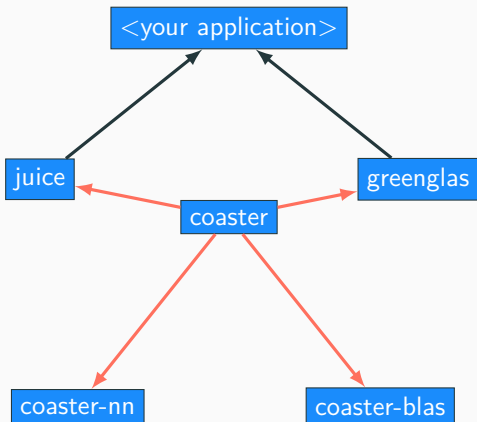
Crates



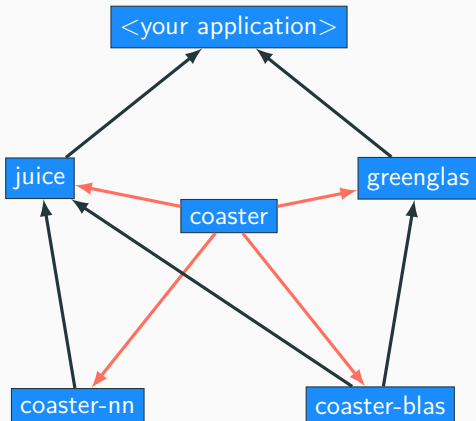
Crates

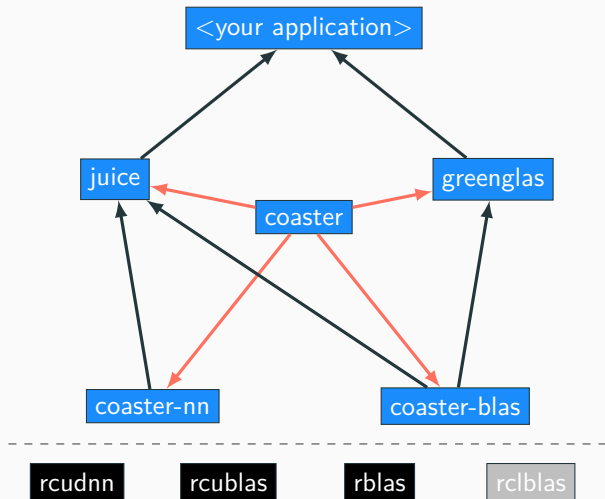


Crates



Crates





A close-up photograph of a mechanical engine component, likely a spark plug or a similar part, with a chain or belt mechanism. The image is highly detailed, showing various metal parts, bolts, and a chain. The lighting is dramatic, with strong highlights and deep shadows, creating a sense of depth and texture. The background is blurred, focusing attention on the foreground components.

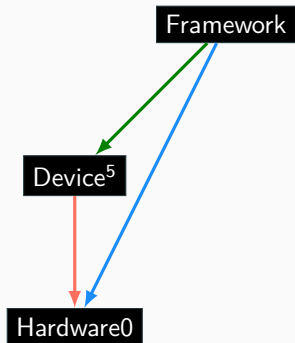
Coaster

Abstraction

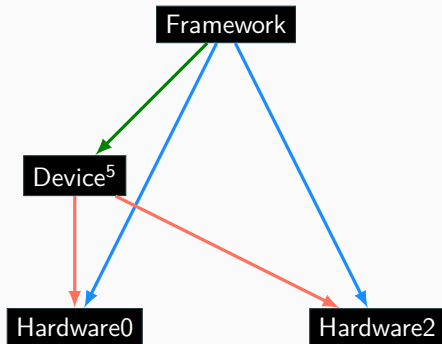
- device memory
- internal structure
- backend specific

```
pub struct SharedTensor<T> {  
    desc: TensorDesc,  
    locations: RefCell<Vec<TensorLocation>>,  
    up_to_date: Cell<BitMap>,  
}
```

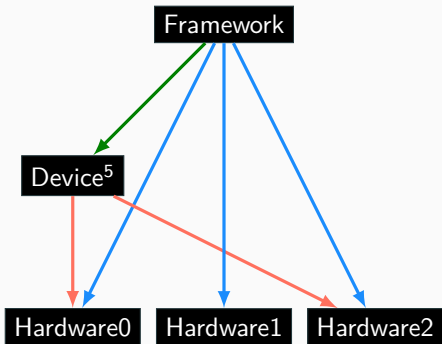
Backend Abstraction⁶



Backend Abstraction⁶

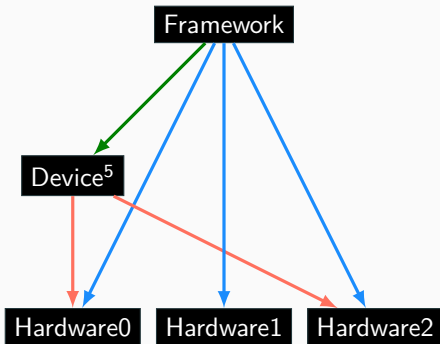


Backend Abstraction⁶

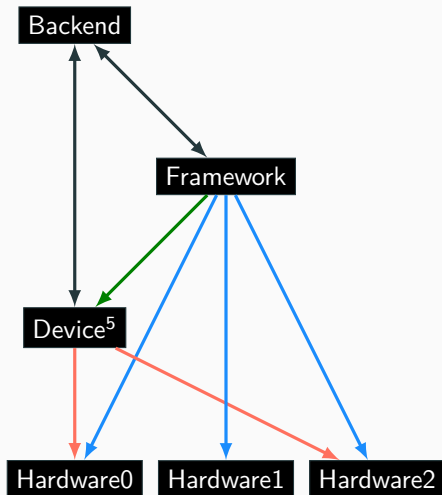


Backend Abstraction⁶

Backend



Backend Abstraction⁶



⁵better HardwareGroup/DeviceGroup

⁶the members, associated types differ!


```
/// Defines a Framework.
pub trait IFramework {
    /// The Hardware representation for this Framework.
    type H: IHardware;
    /// The Device representation for this Framework.
    type D: IDevice + Clone;
    /// The Binary representation for this Framework.
    type B: IBinary + Clone;
```

⁷or better: CPU

```
/// Defines a Framework.  
pub trait IFramework {  
    /// The Hardware representation for this Framework.  
    type H: IHardware;  
    /// The Device representation for this Framework.  
    type D: IDevice + Clone;  
    /// The Binary representation for this Framework.  
    type B: IBinary + Clone;
```

- cuda
- OpenCL
- Native⁷

⁷or better: CPU

```
pub struct Backend<F: IFramework> {  
    framework: Box<F>,  
    device: F::D,  
}
```

```
pub struct Backend<F: IFramework> {  
    framework: Box<F>,  
    device: F::D,  
}
```

Backend is tied to a Framework

Backend has a Device

Device has Hardwares



Juice

**Memory objects with arbitrary (but defined!)
structure and dimensions**

=

Blobs

=

SharedTensor

parameters
=
weights (+ bias)

Architecture

`trait ILayer` common interface that is expected to be impl'd by all layers

⁸mostly a helper to make LayerConfig impl easier

⁹or any other for that matter, just i.e.

Architecture

`trait ILayer` common interface that is expected to be impl'd by all layers

`LayerConfig` handles input and output tracking, name based

⁸mostly a helper to make `LayerConfig` impl easier

⁹or any other for that matter, just i.e.

Architecture

trait ILayer common interface that is expected to be impl'd by all layers

LayerConfig handles input and output tracking, name based

LayerType holds configuration data for the particular layer⁸

⁸mostly a helper to make LayerConfig impl easier

⁹or any other for that matter, just i.e.

Architecture

trait ILayer common interface that is expected to be impl'd by all layers

LayerConfig handles input and output tracking, name based

LayerType holds configuration data for the particular layer⁸

PoolingConfig⁹ holds specific information

⁸mostly a helper to make LayerConfig impl easier

⁹or any other for that matter, just i.e.

Layers implemented

- activation ReLU
- activation Sigmoid
- activation TanH
- common Dropout¹⁰
- common Softmax
- common LogSoftmax
- common Convolution ¹¹
- common Linear
- common Pooling
- util Reshape
- util Flatten

¹⁰unmerged, but it is there pr#13

¹¹the native impl has only forward impl

Containers

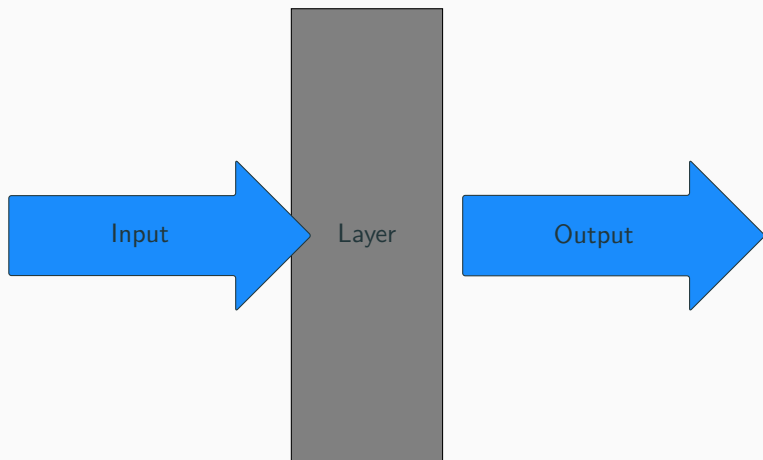
- most common: **Sequential / SequentialConfig**
- **stackable, container in a container in a container** ¹²
- are **Layers/LayerConfigs** too

¹²no docker involved

ILayer

ILayer
=
ComputeOutput
+ ComputeInputGradient
+ ComputeParametersGradient

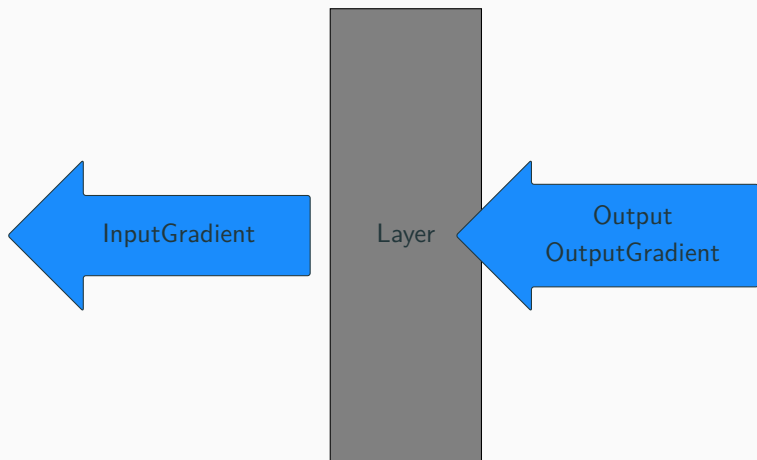
DataFlow: ComputeOutput



Layer: ComputeOutput

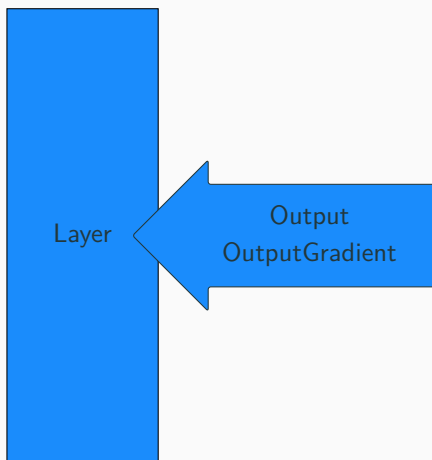
```
/// A Layer that can compute the output for a given input.
pub trait ComputeOutput<T, B: IBackend> {
    /// Compute output for given input
    /// and write them into `output_data`.
    fn compute_output(&self,
        backend: &B,
        weights_data: &[&SharedTensor<T>],
        input_data: &[&SharedTensor<T>],
        output_data: &mut [&mut SharedTensor<T>]);
}
```


DataFlow: ComputeInputGradient



Layer: ComputeInputGradient

```
/// A Layer that can compute the gradient with respect to its input.
pub trait ComputeInputGradient<T, B: IBackend> {
    /// Compute gradients with respect to the inputs
    /// and write them into `input_gradients`.
    fn compute_input_gradient(&self,
                               backend: &B,
                               weights_data: &[&SharedTensor<T>],
                               output_data: &[&SharedTensor<T>],
                               output_gradients: &[&SharedTensor<T>],
                               input_data: &[&SharedTensor<T>],
                               input_gradients: &mut [&mut SharedTensor<T>])
}
```



Layer: ComputeParametersGradient

```
/// A Layer that can compute the gradient with respect to its parameter
pub trait ComputeParametersGradient<T, B: IBackend> {
    /// Compute gradients with respect to the parameters
    /// and write them into `parameters_gradients`.
    fn compute_parameters_gradient(&self,
                                   backend: &B,
                                   output_data: &[&SharedTensor<T>],
                                   output_gradients: &[&SharedTensor<T>],
                                   input_data: &[&SharedTensor<T>],
                                   parameters_gradients: &mut [&mut SharedTensor<T>],
                                   // not required for all layers,
                                   // only those with internal state to be optimized
    )
}
```



Something is still missing.

Something is still missing (?)




Solver

- connects forward and backward layer
- optimized towards objective¹³
- contains almost all tuneable hyperparameters¹⁴

¹³implicitly defined loss function

¹⁴Parameters which are not optimized towards


```
pub struct SolverConfig {  
  pub name: String,  
  pub network: LayerConfig,  
  pub objective: LayerConfig,  
  pub solver: SolverKind,  
  // more..  
}
```

A roller coaster track with two cars is shown against a cloudy sky. The track is dark grey and features several loops and drops. The cars are white with red accents and are positioned on the track. The background is a light blue sky with soft, white clouds. A dark green rectangular box with white text is overlaid on the center of the image.

Be Fast, Stay Fast

- 95% spent in optimized libraries (i.e. cudnn)
- avoid implementing things manually (blas)
- skip operations if possible
- provide in place operations
- never re-allocate device memory if possible / hold device mem as long as possible

Alexnet example (network cfg)

```
extern crate env_logger;
extern crate coaster as co;
extern crate juice;

use co::prelude::*;

let mut cfg = SequentialConfig::default();
cfg.add_input("data", &[128, channels, px_dim, px_dim]);

let conv1_layer_cfg = ConvolutionConfig {
    num_output: 64,
    filter_shape: vec![11],
    padding: vec![2],
    stride: vec![4],
};

cfg.add_layer(LayerConfig::new("conv1", conv1_layer_cfg));
cfg.add_layer(LayerConfig::new("conv1/relu", LayerType::ReLU));
```

Alexnet example (network cfg)

```
cfg.add_layer(LayerConfig::new("conv2",
    ConvolutionConfig {
        num_output: 192,
        filter_shape: vec![5],
        padding: vec![2],
        stride: vec![1],
    }));
cfg.add_layer(LayerConfig::new("conv2/relu",
    LayerType::ReLU));
cfg.add_layer(LayerConfig::new("pool2",
    PoolingConfig {
        mode: PoolingMode::Max,
        filter_shape: vec![3],
        stride: vec![2],
        padding: vec![0],
    }));
```

Alexnet example (network cfg)

```
// more layers ..
cfg.add_layer(LayerConfig::new("fc1",
    LinearConfig { output_size: 4096 }));
cfg.add_layer(LayerConfig::new("fc2",
    LinearConfig { output_size: 4096 }));
cfg.add_layer(LayerConfig::new("fc3",
    LinearConfig { output_size: 1000 }));

// create pseudo probabilities + log
cfg.add_layer(LayerConfig::new("log_softmax", LayerType::LogSoftmax));

// set up backends
let backend = Rc::new(Backend::<Cuda>::default().unwrap());

let mut network = Layer::from_config(backend.clone(),
    &LayerConfig::new("alexnet",
        LayerType::Sequential(cfg)));
```

Alexnet example (classifier cfg)

```
let mut classifier_cfg = SequentialConfig::default();
classifier_cfg.add_input("network_out", &[batch_size, 10]);
classifier_cfg.add_input("label", &[batch_size, 1]);
// set up nll loss
let nll_layer_cfg = NegativeLogLikelihoodConfig { num_classes: 10 };
let nll_cfg = LayerConfig::new("nll",
                               LayerType::NegativeLogLikelihood(nll_layer_cfg));
classifier_cfg.add_layer(nll_cfg);
```

Alexnet example (training/backward pass, 1)

```
// set up solver
let mut solver_cfg = SolverConfig {
    minibatch_size: batch_size,
    base_lr: learning_rate,
    momentum: momentum,
    ..SolverConfig::default()
};
solver_cfg.network = LayerConfig::new("network", cfg);
solver_cfg.objective = LayerConfig::new("classifier",
    classifier_cfg);
let mut solver = Solver::from_config(backend.clone(),
    &solver_cfg);

// set up confusion matrix
let mut confusion = ConfusionMatrix::new(10);
confusion.set_capacity(Some(1000));
```


Alexnet example (training/backward pass, 2)

```
let inp = SharedTensor::<f32>::new(&[batch_size, channels, px_dim, py_dim]);
let label = SharedTensor::<f32>::new(&[batch_size, 1]);
let inp_lock = Arc::new(RwLock::new(inp));
let label_lock = Arc::new(RwLock::new(label));

for _ in 0..(example_count / batch_size) {
    let mut targets = Vec::new();
    for (batch_n, (label_val, input)) in decoded_images
        .by_ref()
        .take(batch_size)
        .enumerate() {
        let mut inp = inp_lock.write().unwrap();
        let mut label = label_lock.write().unwrap();
        write_batch_sample(&mut inp, &input, batch_n);
        write_batch_sample(&mut label, &[label_val], batch_n);
        targets.push(label_val as usize);
    }
}
// ...
```

Alexnet example (training/backward pass, 2)

```
for _ in 0..(example_count / batch_size) {
  // ...
  // train the network!
  let inferred_out = solver.train_minibatch(
    inp_lock.clone(),
    label_lock.clone());
  let mut inferred = inferred_out.write().unwrap();
  let predictions = confusion.get_predictions(&mut inferred);

  confusion.add_samples(&predictions, &targets);
  println!(
    "Last sample: {} | Accuracy {}",
    confusion.samples().iter().last().unwrap(),
    confusion.accuracy()
  );
}
```

Alexnet example (prediction/inference/forward pass)

```
let inp = SharedTensor::new(&[1, channels, px_dim, px_dim]);  
let inp_lock = Arc::new(RwLock::new(inp));  
let predictions = network.forward(&[inp_lock.clone()]);  
// do something with those predictions
```

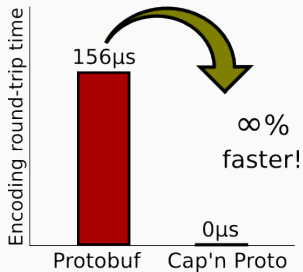
mnist, fashion-mnist


`github.com/spearow/juice-examples`



save/load/store/export

- interface to other languages
- load, store weights and the layer architecture
- it executes *really fast!*





What happened in 2017?

Summary 2017

- Learning a whole lot more about rust ¹⁵
- Merging remaining large pull request ¹⁶
- impl Dropout in coaster-nn for backends cuda and native
- Setting up CI ci.spearow.io / @sirmergealot 🇩🇪
- Auto generating config files for CI spearow/ci
- Fight a lot with CI, get GPU hardware access from containers from garden hypervisor ¹⁷

¹⁵build.rs, bindgen, PhantomData, Any trait and more

¹⁶hands down @alexandermorozov

¹⁷Still an ugly hack



Roadmap 2018

- Recursive Neural Networks / Long Term Short Term Memory [coaster-nn]
- OpenCL Backend [juice/coaster-nn]
- Autodiff [juice/coaster-nn]
- Accuracy enum instead of type usage [juice/greenglas/coaster]

- Honour biases [juice]
- Gradient calculation for ND-convolution native backend [coaster-nn]
- Regression examples [juice-examples]
- Add preprocessing filters ¹⁸ [greenglas]

¹⁸FFT, wavelet, denoise, addnoise, cutout, rescale, blur, ..

Questions?

bernhard@sparrow.io

sparrow.io

gitter.im/sparrow/juice

Credits

Presentation by Bernhard Schuster - ahoi.io
Theme by Matze Vogelsang and contributors -
<https://bloerg.net>

The theme *itself* is licensed under a Creative Commons
Attribution-ShareAlike 4.0 International License.

The images are taken from *unsplash* and are free to use for whatever you
want.

Font Awesome icons SIL OFL 1.1

Icons displayed made by Bernhard Schuster.

