

# Overview

Welcome to The Burn Book 🙌

This book will help you get started with the Burn deep learning framework, whether you are an advanced user or a beginner. We have crafted some sections for you:

- [Basic Workflow: From Training to Inference](#): We'll start with the fundamentals, guiding you through the entire workflow, from training your models to deploying them for inference. This section lays the groundwork for your Burn expertise.
- [Building Blocks](#): Dive deeper into Burn's core components, understanding how they fit together. This knowledge forms the basis for more advanced usage and customization.
- [Saving & Loading Models](#): Learn how to easily save and load your trained models.
- [Custom Training Loop](#): Gain the power to customize your training loops, fine-tuning your models to meet your specific requirements. This section empowers you to harness Burn's flexibility to its fullest.
- [Importing Models](#): Learn how to import ONNX and PyTorch models, expanding your compatibility with other deep learning ecosystems.
- [Advanced](#): Finally, venture into advanced topics, exploring Burn's capabilities at their peak. This section caters to those who want to push the boundaries of what's possible with Burn.

Throughout the book, we assume a basic understanding of deep learning concepts, but we may refer to additional material when it seems appropriate.

# Why Burn?

Why bother with the effort of creating an entirely new deep learning framework from scratch when PyTorch, TensorFlow, and other frameworks already exist? Spoiler alert: **Burn isn't merely a replication of PyTorch or TensorFlow in Rust.** It represents a novel approach, placing significant emphasis on making the right compromises in the right areas to facilitate exceptional flexibility, high performance, and a seamless developer experience. **Burn isn't a framework specialized for only one type of application, it is designed to serve as a versatile framework suitable for a wide range of research and production uses.** The foundation of Burn's design revolves around three key user profiles:

**Machine Learning Researchers** require tools to construct and execute experiments efficiently. It's essential for them to iterate quickly on their ideas and design testable experiments which can help them discover new findings. The framework should facilitate the swift implementation of cutting-edge research while ensuring fast execution for testing.

**Machine Learning Engineers** are another important demographic to keep in mind. Their focus leans less on swift implementation and more on establishing robustness, seamless deployment, and cost-effective operations. They seek dependable, economical models capable of achieving objectives without excessive expense. The whole machine learning workflow—from training to inference—must be as efficient as possible with minimal unpredictable behavior.

**Low level Software Engineers** working with hardware vendors want their processing units to run models as fast as possible to gain competitive advantage. This endeavor involves harnessing hardware-specific features such as Tensor Core for Nvidia. Since they are mostly working at a system level, they want to have absolute control over how the computation will be executed.

The goal of Burn is to satisfy all of those personas!

# Getting Started

Burn is a deep learning framework in the Rust programming language. Therefore, it goes without saying that one must understand the basic notions of Rust. Reading the first chapters of the [Rust Book](#) is recommended, but don't worry if you're just starting out. We'll try to provide as much context and reference to external resources when required. Just look out for the 🦀 **Rust Note** indicators.

## Installing Rust

For installation instructions, please refer to the [installation page](#). It explains in details the most convenient way for you to install Rust on your computer, which is the very first thing to do to start using Burn.

## Creating a Burn application

Once Rust is correctly installed, create a new Rust application by using Rust's build system and package manager Cargo. It is automatically installed with Rust.

### ► 🦀 Cargo Cheat Sheet

In the directory of your choice, run the following:

```
cargo new my_burn_app
```

This will initialize the `my_burn_app` project directory with a `Cargo.toml` file and a `src` directory with an auto-generated `main.rs` file inside. Head inside the directory to check:

```
cd my_burn_app
```

Then, add Burn as a dependency:

```
cargo add burn --features wgpu
```

Finally, compile the local package by executing the following:

```
cargo build
```

That's it, you're ready to start! You have a project configured with Burn and the WGPU backend, which allows to execute low-level operations on any platform using the GPU.

## Writing a code snippet

The `src/main.rs` was automatically generated by Cargo, so let's replace its content with the following:

```
use burn::tensor::Tensor;
use burn::backend::Wgpu;

// Type alias for the backend to use.
type Backend = Wgpu;

fn main() {
    let device = Default::default();
    // Creation of two tensors, the first with explicit values and the second
    one with ones, with the same shape as the first
    let tensor_1 = Tensor::<Backend, 2>::from_data([[2., 3.], [4., 5.]],
&device);
    let tensor_2 = Tensor::<Backend, 2>::ones_like(&tensor_1);

    // Print the element-wise addition (done with the WGPU backend) of the two
    tensors.
    println!("{}", tensor_1 + tensor_2);
}
```

### ► 🦀 Use Declarations

### ► 🦀 Generic Data Types

By running `cargo run`, you should now see the result of the addition:

```
Tensor {
  data:
  [[3.0, 4.0],
   [5.0, 6.0]],
  shape: [2, 2],
  device: BestAvailable,
  backend: "wgpu",
  kind: "Float",
  dtype: "f32",
}
```

While the previous example is somewhat trivial, the upcoming basic workflow section will walk you through a much more relevant example for deep learning applications.

## Using prelude

Burn comes with a variety of things in its core library. When creating a new model or using an existing one for inference, you may need to import every single component you used, which could be a little verbose.

To address it, a `prelude` module is provided, allowing you to easily import commonly used structs and macros as a group:

```
use burn::prelude::*;
```

which is equal to:

```
use burn::{
  config::Config,
  module::Module,
  nn,
  tensor::{
    backend::Backend, Bool, Device, ElementConversion, Float, Int, Shape,
    Tensor,
    TensorData,
  },
};
```

For the sake of simplicity, the subsequent chapters of this book will all use this form of importing except in the [Building Blocks](#) chapter, as explicit importing aids users in grasping the usage of particular structures and macros.

# Examples

In the [next chapter](#) you'll have the opportunity to implement the whole Burn `guide` example yourself in a step by step manner.

Many additional Burn examples are available in the [examples](#) directory. Burn examples are organized as library crates with one or more examples that are executable binaries. An example can then be executed using the following cargo command line in the root of the Burn repository:

```
cargo run --example <example name>
```

To learn more about crates and examples, read the Rust section below.

## ► About Rust crates

The following additional examples are currently available if you want to check them out:

Example	Description
<a href="#">Custom CSV Dataset</a>	Implements a dataset to parse CSV data for a regression task.
<a href="#">Regression</a>	Trains a simple MLP on the California Housing dataset to predict the median house value for a district.
<a href="#">Custom Image Dataset</a>	Trains a simple CNN on custom image dataset following a simple folder structure.
<a href="#">Custom Renderer</a>	Implements a custom renderer to display the <a href="#">Learner</a> progress.
<a href="#">Image Classification Web</a>	Image classification web browser demo using Burn, WGPU and WebAssembly.
<a href="#">MNIST Inference on Web</a>	An interactive MNIST inference demo in the browser. The demo is available <a href="#">online</a> .
<a href="#">MNIST Training</a>	Demonstrates how to train a custom <a href="#">Module</a> (MLP) with the <a href="#">Learner</a> configured to log metrics and keep training checkpoints.
<a href="#">Named Tensor</a>	Performs operations with the experimental <code>NamedTensor</code> feature.

Example	Description
<a href="#">ONNX Import Inference</a>	Imports an ONNX model pre-trained on MNIST to perform inference on a sample image with Burn.
<a href="#">PyTorch Import Inference</a>	Imports a PyTorch model pre-trained on MNIST to perform inference on a sample image with Burn.
<a href="#">Text Classification</a>	Trains a text classification transformer model on the AG News or DbPedia datasets. The trained model can then be used to classify a text sample.
<a href="#">Text Generation</a>	Trains a text generation transformer model on the DbPedia dataset.

For more information on each example, see their respective `README.md` file. Be sure to check out the [examples](#) directory for an up-to-date list.

Note that some examples use the [datasets library by HuggingFace](#) to download the datasets required in the examples. This is a Python library, which means that you will need to install Python before running these examples. This requirement will be clearly indicated in the example's README when applicable.

# Guide

This guide will walk you through the process of creating a custom model built with Burn. We will train a simple convolutional neural network model on the MNIST dataset and prepare it for inference.

For clarity, we sometimes omit imports in our code snippets. For more details, please refer to the corresponding code in the `examples/guide` [directory](#). We reproduce this example in a step-by-step fashion, from dataset creation to modeling and training in the following sections. It is recommended to use the capabilities of your IDE or text editor to automatically add the missing imports as you add the code snippets to your code.

Be sure to checkout the git branch corresponding to the version of Burn you are using to follow this guide.

The current version of Burn is `0.15` and the corresponding branch to checkout is [release/0.15](#).

The code for this demo can be executed from Burn's base directory using the command:

```
cargo run --example guide
```

## Key Learnings

- Creating a project
- Creating neural network models
- Importing and preparing datasets
- Training models on data
- Choosing a backend
- Using a model for inference



# Model

The first step is to create a project and add the different Burn dependencies. Start by creating a new project with Cargo:

```
cargo new guide
```

As [mentioned previously](#), this will initialize your `guide` project directory with a `Cargo.toml` and a `src/main.rs` file.

In the `Cargo.toml` file, add the `burn` dependency with `train`, `wgpu` and `vision` features. Then run `cargo build` to build the project and import all the dependencies.

```
[package]
name = "guide"
version = "0.1.0"
edition = "2021"

[dependencies]
burn = { version = "~0.15", features = ["train", "wgpu", "vision"] }
```

Our goal will be to create a basic convolutional neural network used for image classification. We will keep the model simple by using two convolution layers followed by two linear layers, some pooling and ReLU activations. We will also use dropout to improve training performance.

Let us start by defining our model struct in a new file `src/model.rs`.

```

use burn::{
    nn::{
        conv::{Conv2d, Conv2dConfig},
        pool::{AdaptiveAvgPool2d, AdaptiveAvgPool2dConfig},
        Dropout, DropoutConfig, Linear, LinearConfig, Relu,
    },
    prelude::*,
};

#[derive(Module, Debug)]
pub struct Model<B: Backend> {
    conv1: Conv2d<B>,
    conv2: Conv2d<B>,
    pool: AdaptiveAvgPool2d,
    dropout: Dropout,
    linear1: Linear<B>,
    linear2: Linear<B>,
    activation: Relu,
}

```

There are two major things going on in this code sample.

1. You can create a deep learning module with the `#[derive(Module)]` attribute on top of a struct. This will generate the necessary code so that the struct implements the `Module` trait. This trait will make your module both trainable and (de)serializable while adding related functionalities. Like other attributes often used in Rust, such as `Clone`, `PartialEq` or `Debug`, each field within the struct must also implement the `Module` trait.

► 🦀 **Trait**

► 🦀 **Derive Macro**

2. Note that the struct is generic over the `Backend` trait. The backend trait abstracts the underlying low level implementations of tensor operations, allowing your new model to run on any backend. Contrary to other frameworks, the backend abstraction isn't determined by a compilation flag or a device type. This is important because you can extend the functionalities of a specific backend (see [backend extension section](#)), and it allows for an innovative [autodiff system](#). You can also change backend during runtime, for instance to compute training metrics on a cpu backend while using a gpu one only to train the model. In our example, the backend in use will be determined later on.

► 🦀 **Trait Bounds**

Note that each time you create a new file in the `src` directory you also need to explicitly add this module to the `main.rs` file. For instance after creating the `model.rs`, you need to add the following at the top of the main file:

```
mod model;
```

Next, we need to instantiate the model for training.

```
#[derive(Config, Debug)]
pub struct ModelConfig {
    num_classes: usize,
    hidden_size: usize,
    #[config(default = "0.5")]
    dropout: f64,
}

impl ModelConfig {
    /// Returns the initialized model.
    pub fn init<B: Backend>(&self, device: &B::Device) -> Model<B> {
        Model {
            conv1: Conv2dConfig::new([1, 8], [3, 3]).init(device),
            conv2: Conv2dConfig::new([8, 16], [3, 3]).init(device),
            pool: AdaptiveAvgPool2dConfig::new([8, 8]).init(),
            activation: Relu::new(),
            linear1: LinearConfig::new(16 * 8 * 8,
self.hidden_size).init(device),
            linear2: LinearConfig::new(self.hidden_size,
self.num_classes).init(device),
            dropout: DropoutConfig::new(self.dropout).init(),
        }
    }
}
```

At a glance, you can view the model configuration by printing the model instance:

```
use crate::model::ModelConfig;
use burn::backend::Wgpu;

fn main() {
    type MyBackend = Wgpu<f32, i32>;

    let device = Default::default();
    let model = ModelConfig::new(10, 512).init::<MyBackend>(&device);

    println!("{}", model);
}
```

Output:

```
Model {  
  conv1: Conv2d {stride: [1, 1], kernel_size: [3, 3], dilation: [1, 1], groups:  
1, padding: Valid, params: 80}  
  conv2: Conv2d {stride: [1, 1], kernel_size: [3, 3], dilation: [1, 1], groups:  
1, padding: Valid, params: 1168}  
  pool: AdaptiveAvgPool2d {output_size: [8, 8]}  
  dropout: Dropout {prob: 0.5}  
  linear1: Linear {d_input: 1024, d_output: 512, bias: true, params: 524800}  
  linear2: Linear {d_input: 512, d_output: 10, bias: true, params: 5130}  
  activation: Relu  
  params: 531178  
}
```

## ► References

When creating a custom neural network module, it is often a good idea to create a config alongside the model struct. This allows you to define default values for your network, thanks to the `Config` attribute. The benefit of this attribute is that it makes the configuration serializable, enabling you to painlessly save your model hyperparameters, enhancing your experimentation process. Note that a constructor will automatically be generated for your configuration, which will take in as input values the parameters which do not have default values: `let config = ModelConfig::new(num_classes, hidden_size);`. The default values can be overridden easily with builder-like methods: (e.g `config.with_dropout(0.2);`)

The first implementation block is related to the initialization method. As we can see, all fields are set using the configuration of the corresponding neural network's underlying module. In this specific case, we have chosen to expand the tensor channels from 1 to 8 with the first layer, then from 8 to 16 with the second layer, using a kernel size of 3 on all dimensions. We also use the adaptive average pooling module to reduce the dimensionality of the images to an 8 by 8 matrix, which we will flatten in the forward pass to have a 1024 (16 \_ 8 \_ 8) resulting tensor.

Now let's see how the forward pass is defined.

```

impl<B: Backend> Model<B> {
  /// # Shapes
  ///   - Images [batch_size, height, width]
  ///   - Output [batch_size, num_classes]
  pub fn forward(&self, images: Tensor<B, 3>) -> Tensor<B, 2> {
    let [batch_size, height, width] = images.dims();

    // Create a channel at the second dimension.
    let x = images.reshape([batch_size, 1, height, width]);

    let x = self.conv1.forward(x); // [batch_size, 8, _, _]
    let x = self.dropout.forward(x);
    let x = self.conv2.forward(x); // [batch_size, 16, _, _]
    let x = self.dropout.forward(x);
    let x = self.activation.forward(x);

    let x = self.pool.forward(x); // [batch_size, 16, 8, 8]
    let x = x.reshape([batch_size, 16 * 8 * 8]);
    let x = self.linear1.forward(x);
    let x = self.dropout.forward(x);
    let x = self.activation.forward(x);

    self.linear2.forward(x) // [batch_size, num_classes]
  }
}

```

For former PyTorch users, this might feel very intuitive, as each module is directly incorporated into the code using an eager API. Note that no abstraction is imposed for the forward method. You are free to define multiple forward functions with the names of your liking. Most of the neural network modules already built with Burn use the `forward` nomenclature, simply because it is standard in the field.

Similar to neural network modules, the `Tensor` struct given as a parameter also takes the Backend trait as a generic argument, alongside its dimensionality. Even if it is not used in this specific example, it is possible to add the kind of the tensor as a third generic argument. For example, a 3-dimensional Tensor of different data types(float, int, bool) would be defined as following:

```

Tensor<B, 3> // Float tensor (default)
Tensor<B, 3, Float> // Float tensor (explicit)
Tensor<B, 3, Int> // Int tensor
Tensor<B, 3, Bool> // Bool tensor

```

Note that the specific element type, such as `f16`, `f32` and the likes, will be defined later with the backend.

# Data

Typically, one trains a model on some dataset. Burn provides a library of very useful dataset sources and transformations, such as Hugging Face dataset utilities that allow to download and store data into an SQLite database for extremely efficient data streaming and storage.

For this guide though, we will use the MNIST dataset from `burn::data::dataset::vision` which requires no external dependency.

To iterate over a dataset efficiently, we will define a struct which will implement the `Batcher` trait. The goal of a batcher is to map individual dataset items into a batched tensor that can be used as input to our previously defined model.

Let us start by defining our dataset functionalities in a file `src/data.rs`. We shall omit some of the imports for brevity, but the full code for following this guide can be found at [examples/guide/ directory](#).

```
use burn::{
    data::{dataloader::batcher::Batcher, dataset::vision::MnistItem},
    prelude::*,
};

#[derive(Clone)]
pub struct MnistBatcher<B: Backend> {
    device: B::Device,
}

impl<B: Backend> MnistBatcher<B> {
    pub fn new(device: B::Device) -> Self {
        Self { device }
    }
}
```

This codeblock defines a batcher struct with the device in which the tensor should be sent before being passed to the model. Note that the device is an associative type of the `Backend` trait since not all backends expose the same devices. As an example, the Libtorch-based backend exposes `Cuda(gpu_index)`, `Cpu`, `Vulkan` and `Metal` devices, while the `ndarray` backend only exposes the `Cpu` device.

Next, we need to actually implement the batching logic.

```

#[derive(Clone, Debug)]
pub struct MnistBatch<B: Backend> {
    pub images: Tensor<B, 3>,
    pub targets: Tensor<B, 1, Int>,
}

impl<B: Backend> Batcher<MnistItem, MnistBatch<B>> for MnistBatcher<B> {
    fn batch(&self, items: Vec<MnistItem>) -> MnistBatch<B> {
        let images = items
            .iter()
            .map(|item| TensorData::from(item.image).convert:::<B::FloatElem>())
            .map(|data| Tensor:::<B, 2>::from_data(data, &self.device))
            .map(|tensor| tensor.reshape([1, 28, 28]))
            // Normalize: make between [0,1] and make the mean=0 and std=1
            // values mean=0.1307,std=0.3081 are from the PyTorch MNIST example
            // https://github.com/pytorch/examples/
            blob/54f4572509891883a947411fd7239237dd2a39c3/mnist/main.py#L122
            .map(|tensor| ((tensor / 255) - 0.1307) / 0.3081)
            .collect();

        let targets = items
            .iter()
            .map(|item| {
                Tensor:::<B, 1, Int>::from_data(
                    [(item.label as i64).elem:::<B::IntElem>()],
                    &self.device,
                )
            })
            .collect();

        let images = Tensor:::cat(images, 0).to_device(&self.device);
        let targets = Tensor:::cat(targets, 0).to_device(&self.device);

        MnistBatch { images, targets }
    }
}

```

## ► Iterators and Closures

In the previous example, we implement the `Batcher` trait with a list of `MnistItem` as input and a single `MnistBatch` as output. The batch contains the images in the form of a 3D tensor, along with a targets tensor that contains the indexes of the correct digit class. The first step is to parse the image array into a `TensorData` struct. Burn provides the `TensorData` struct to encapsulate tensor storage information without being specific for a backend. When creating a tensor from data, we often need to convert the data precision to the current backend in use. This can be done with the `.convert()` method (in this example, the data is converted backend's float element type `B::FloatElem`). While importing the

`burn::tensor::ElementConversion` trait, you can call `.elem()` on a specific number to convert it to the current backend element type in use.



# Training

We are now ready to write the necessary code to train our model on the MNIST dataset. We shall define the code for this training section in the file: `src/training.rs`.

Instead of a simple tensor, the model should output an item that can be understood by the learner, a struct whose responsibility is to apply an optimizer to the model. The output struct is used for all metrics calculated during the training. Therefore it should include all the necessary information to calculate any metric that you want for a task.

Burn provides two basic output types: `ClassificationOutput` and `RegressionOutput`. They implement the necessary trait to be used with metrics. It is possible to create your own item, but it is beyond the scope of this guide.

Since the MNIST task is a classification problem, we will use the `ClassificationOutput` type.

```
impl<B: Backend> Model<B> {
    pub fn forward_classification(
        &self,
        images: Tensor<B, 3>,
        targets: Tensor<B, 1, Int>,
    ) -> ClassificationOutput<B> {
        let output = self.forward(images);
        let loss = CrossEntropyLossConfig::new()
            .init(&output.device())
            .forward(output.clone(), targets.clone());

        ClassificationOutput::new(loss, output, targets)
    }
}
```

As evident from the preceding code block, we employ the cross-entropy loss module for loss calculation, without the inclusion of any padding token. We then return the classification output containing the loss, the output tensor with all logits and the targets.

Please take note that tensor operations receive owned tensors as input. For reusing a tensor multiple times, you need to use the `clone()` function. There's no need to worry; this process won't involve actual copying of the tensor data. Instead, it will simply indicate that the tensor is employed in multiple instances, implying that certain operations won't be performed in place. In summary, our API has been designed with owned tensors to optimize performance.

Moving forward, we will proceed with the implementation of both the training and validation

steps for our model.

```
impl<B: AutodiffBackend> TrainStep<MnistBatch<B>, ClassificationOutput<B>> for
Model<B> {
  fn step(&self, batch: MnistBatch<B>) ->
TrainOutput<ClassificationOutput<B>> {
    let item = self.forward_classification(batch.images, batch.targets);

    TrainOutput::new(self, item.loss.backward(), item)
  }
}

impl<B: Backend> ValidStep<MnistBatch<B>, ClassificationOutput<B>> for Model<B>
{
  fn step(&self, batch: MnistBatch<B>) -> ClassificationOutput<B> {
    self.forward_classification(batch.images, batch.targets)
  }
}
```

Here we define the input and output types as generic arguments in the `TrainStep` and `ValidStep`. We will call them `MnistBatch` and `ClassificationOutput`. In the training step, the computation of gradients is straightforward, necessitating a simple invocation of `backward()` on the loss. Note that contrary to PyTorch, gradients are not stored alongside each tensor parameter, but are rather returned by the backward pass, as such: `let gradients = loss.backward();`. The gradient of a parameter can be obtained with the `grad` function: `let grad = tensor.grad(&gradients);`. Although it is not necessary when using the learner struct and the optimizers, it can prove to be quite useful when debugging or writing custom training loops. One of the differences between the training and the validation steps is that the former requires the backend to implement `AutodiffBackend` and not just `Backend`. Otherwise, the `backward` function is not available, as the backend does not support autodiff. We will see later how to create a backend with autodiff support.

## ► 🦀 Generic Type Constraints in Method Definitions

Let us move on to establishing the practical training configuration.

```
#[derive(Config)]
pub struct TrainingConfig {
    pub model: ModelConfig,
    pub optimizer: AdamConfig,
    #[config(default = 10)]
    pub num_epochs: usize,
    #[config(default = 64)]
    pub batch_size: usize,
    #[config(default = 4)]
    pub num_workers: usize,
    #[config(default = 42)]
    pub seed: u64,
    #[config(default = 1.0e-4)]
    pub learning_rate: f64,
}

fn create_artifact_dir(artifact_dir: &str) {
    // Remove existing artifacts before to get an accurate learner summary
    std::fs::remove_dir_all(artifact_dir).ok();
    std::fs::create_dir_all(artifact_dir).ok();
}

pub fn train<B: AutodiffBackend>(artifact_dir: &str, config: TrainingConfig,
device: B::Device) {
    create_artifact_dir(artifact_dir);
    config
        .save(format!("{artifact_dir}/config.json"))
        .expect("Config should be saved successfully");

    B::seed(config.seed);

    let batcher_train = MnistBatcher::<B>::new(device.clone());
    let batcher_valid = MnistBatcher::<B::InnerBackend>::new(device.clone());

    let dataloader_train = DataLoaderBuilder::new(batcher_train)
        .batch_size(config.batch_size)
        .shuffle(config.seed)
        .num_workers(config.num_workers)
        .build(MnistDataset::train());

    let dataloader_test = DataLoaderBuilder::new(batcher_valid)
        .batch_size(config.batch_size)
        .shuffle(config.seed)
        .num_workers(config.num_workers)
        .build(MnistDataset::test());

    let learner = LearnerBuilder::new(artifact_dir)
        .metric_train_numeric(AccuracyMetric::new())
        .metric_valid_numeric(AccuracyMetric::new())
        .metric_train_numeric(LossMetric::new())
        .metric_valid_numeric(LossMetric::new())
        .with_file_checkpoint(CompactRecorder::new())
```

```

        .devices(vec![device.clone()])
        .num_epochs(config.num_epochs)
        .summary()
        .build(
            config.model.init::<B>(&device),
            config.optimizer.init(),
            config.learning_rate,
        );

let model_trained = learner.fit(dataloader_train, dataloader_test);

model_trained
    .save_file(format!("{artifact_dir}/model"), &CompactRecorder::new())
    .expect("Trained model should be saved successfully");
}

```

It is a good practice to use the `Config` derive to create the experiment configuration. In the `train` function, the first thing we are doing is making sure the `artifact_dir` exists, using the standard rust library for file manipulation. All checkpoints, logging and metrics will be stored under this directory. We then initialize our dataloaders using our previously created batcher. Since no automatic differentiation is needed during the validation phase, the backend used for the corresponding batcher is `B::InnerBackend` (see [Backend](#)). The autodiff capabilities are available through a type system, making it nearly impossible to forget to deactivate gradient calculation.

Next, we create our learner with the accuracy and loss metric on both training and validation steps along with the device and the epoch. We also configure the checkpointer using the `CompactRecorder` to indicate how weights should be stored. This struct implements the `Recorder` trait, which makes it capable of saving records for persistency.

We then build the learner with the model, the optimizer and the learning rate. Notably, the third argument of the `build` function should actually be a learning rate *scheduler*. When provided with a float as in our example, it is automatically transformed into a *constant* learning rate scheduler. The learning rate is not part of the optimizer config as it is often done in other frameworks, but rather passed as a parameter when executing the optimizer step. This avoids having to mutate the state of the optimizer and is therefore more functional. It makes no difference when using the learner struct, but it will be an essential nuance to grasp if you implement your own training loop.

Once the learner is created, we can simply call `fit` and provide the training and validation dataloaders. For the sake of simplicity in this example, we employ the test set as the validation set; however, we do not recommend this practice for actual usage.

Finally, the trained model is returned by the `fit` method. The trained weights are then saved using the `CompactRecorder`. This recorder employs the `MessagePack` format with half

precision, `f16` for floats and `i16` for integers. Other recorders are available, offering support for various formats, such as `BinCode` and `JSON`, with or without compression. Any backend, regardless of precision, can load recorded data of any kind.

# Backend

We have effectively written most of the necessary code to train our model. However, we have not explicitly designated the backend to be used at any point. This will be defined in the main entrypoint of our program, namely the `main` function defined in `src/main.rs`.

```
use crate::{model::ModelConfig, training::TrainingConfig};
use burn::{
    backend::{Autodiff, Wgpu},
    optim::AdamConfig,
};

fn main() {
    type MyBackend = Wgpu<f32, i32>;
    type MyAutodiffBackend = Autodiff<MyBackend>;

    let device = burn::backend::wgpu::WgpuDevice::default();
    let artifact_dir = "/tmp/guide";
    crate::training::train::<MyAutodiffBackend>(
        artifact_dir,
        TrainingConfig::new(ModelConfig::new(10, 512), AdamConfig::new()),
        device.clone(),
    );
}
```

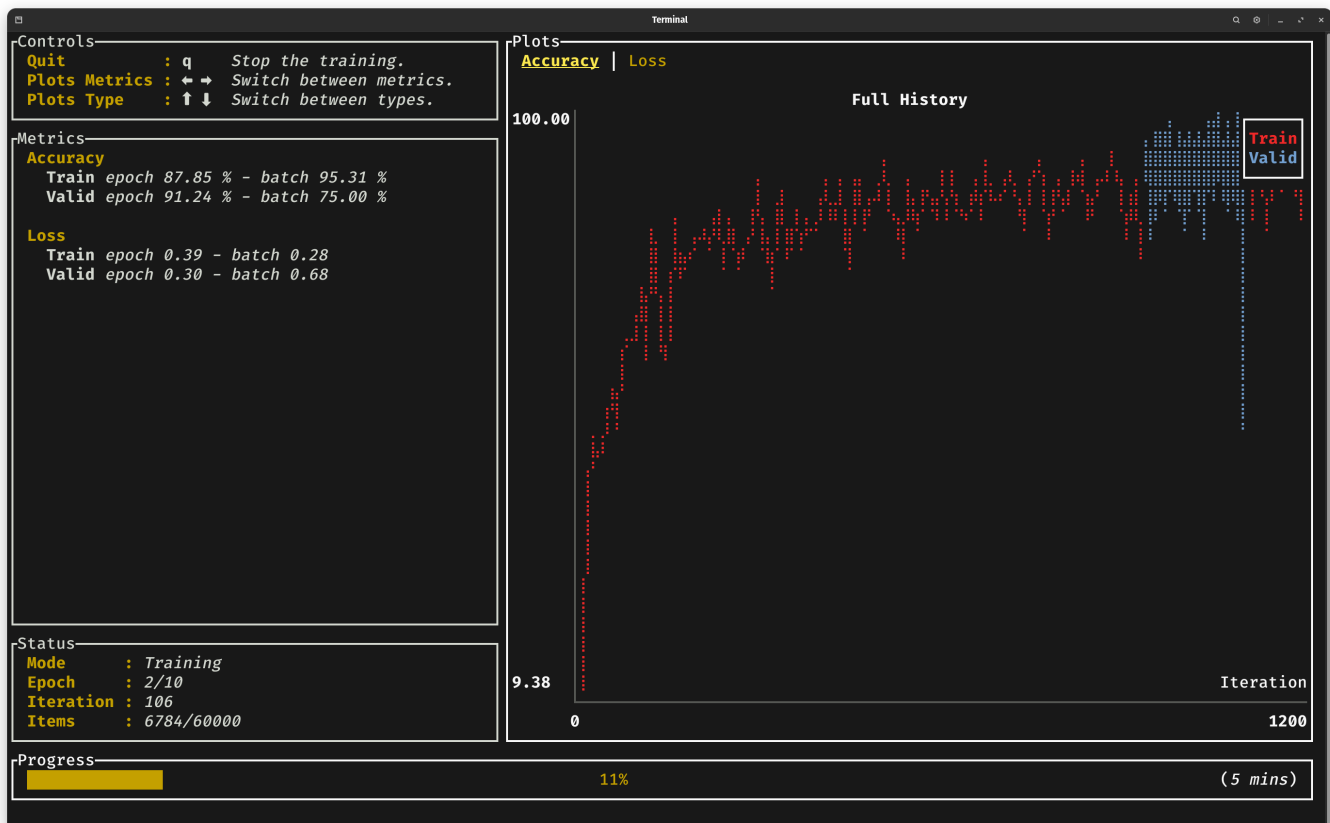
In this code snippet, we use the `wgpu` backend which is compatible with any operating system and will use the GPU. For other options, see the Burn README. This backend type takes the graphics API, the float type and the int type as generic arguments that will be used during the training. The autodiff backend is simply the same backend, wrapped within the `Autodiff` struct which imparts differentiability to any backend.

We call the `train` function defined earlier with a directory for artifacts, the configuration of the model (the number of digit classes is 10 and the hidden dimension is 512), the optimizer configuration which in our case will be the default Adam configuration, and the device which can be obtained from the backend.

You can now train your freshly created model with the command:

```
cargo run --release
```

When running your project with the commande above, you should see the training progression through a basic CLI dashboard:



# Inference

Now that we have trained our model, the next natural step is to use it for inference.

You need two things in order to load weights for a model: the model's record and the model's config. Since parameters in Burn are lazy initialized, no allocation and GPU/CPU kernels are executed by the `ModelConfig::init` function. The weights are initialized when used for the first time, therefore you can safely use `config.init(device).load_record(record)` without any meaningful performance cost. Let's create a simple `infer` method in a new file `src/inference.rs` which we will use to load our trained model.

```
pub fn infer<B: Backend>(artifact_dir: &str, device: B::Device, item:
MnistItem) {
    let config = TrainingConfig::load(format!("{artifact_dir}/config.json"))
        .expect("Config should exist for the model");
    let record = CompactRecorder::new()
        .load(format!("{artifact_dir}/model").into(), &device)
        .expect("Trained model should exist");

    let model = config.model.init::<B>(&device).load_record(record);

    let label = item.label;
    let batcher = MnistBatcher::new(device);
    let batch = batcher.batch(vec![item]);
    let output = model.forward(batch.images);
    let predicted = output.argmax(1).flatten::<1>(0, 1).into_scalar();

    println!("Predicted {} Expected {}", predicted, label);
}
```

The first step is to load the configuration of the training to fetch the correct model configuration. Then we can fetch the record using the same recorder as we used during training. Finally we can init the model with the configuration and the record. For simplicity we can use the same batcher used during the training to pass from a `MnistItem` to a tensor.

By running the `infer` function, you should see the predictions of your model!

Add the call to `infer` to the `main.rs` file after the `train` function call:



```
crate::inference::infer::<MyBackend>(
    artifact_dir,
    device,
    burn::data::dataset::vision::MnistDataset::test()
        .get(42)
        .unwrap(),
);
```

The number 42 is the index of the image in the MNIST dataset. You can explore and verify them using this [MNIST viewer](#).

# Conclusion

In this short guide, we've introduced you to the fundamental building blocks for getting started with Burn. While there's still plenty to explore, our goal has been to provide you with the essential knowledge to kickstart your productivity within the framework.

# Building Blocks

In this section, we'll guide you through the core elements that make up Burn. We'll walk you through the key components that serve as the building blocks of the framework and your future projects.

As you explore Burn, you might notice that we occasionally draw comparisons to PyTorch. We believe it can provide a smoother learning curve and help you grasp the nuances more effectively.

# Backend

Nearly everything in Burn is based on the `Backend` trait, which enables you to run tensor operations using different implementations without having to modify your code. While a backend may not necessarily have autodiff capabilities, the `AutodiffBackend` trait specifies when autodiff is needed. This trait not only abstracts operations but also tensor, device, and element types, providing each backend the flexibility they need. It's worth noting that the trait assumes eager mode since burn fully supports dynamic graphs. However, we may create another API to assist with integrating graph-based backends, without requiring any changes to the user's code.

Users are not expected to directly use the backend trait methods, as it is primarily designed with backend developers in mind rather than Burn users. Therefore, most Burn userland APIs are generic across backends. This approach helps users discover the API more organically with proper autocomplete and documentation.

# Tensor

As previously explained in the [model section](#), the `Tensor` struct has 3 generic arguments: the backend `B`, the dimensionality `D`, and the data type.

```
Tensor<B, D>           // Float tensor (default)
Tensor<B, D, Float>    // Explicit float tensor
Tensor<B, D, Int>      // Int tensor
Tensor<B, D, Bool>     // Bool tensor
```

Note that the specific element types used for `Float`, `Int`, and `Bool` tensors are defined by backend implementations.

Burn Tensors are defined by the number of dimensions `D` in its declaration as opposed to its shape. The actual shape of the tensor is inferred from its initialization. For example, a `Tensor` of size (5,) is initialized as below:

```
let floats = [1.0, 2.0, 3.0, 4.0, 5.0];

// Get the default device
let device = Default::default();

// correct: Tensor is 1-Dimensional with 5 elements
let tensor_1 = Tensor::<Backend, 1>::from_floats(floats, &device);

// incorrect: let tensor_1 = Tensor::<Backend, 5>::from_floats(floats,
&device);
// this will lead to an error and is for creating a 5-D tensor
```

## Initialization

Burn Tensors are primarily initialized using the `from_data()` method which takes the `TensorData` struct as input. The `TensorData` struct has two public fields: `shape` and `dtype`. The `value`, now stored as bytes, is private but can be accessed via any of the following methods: `as_slice`, `as_mut_slice`, `to_vec` and `iter`. To retrieve the data from a tensor, the method `.to_data()` should be employed when intending to reuse the tensor afterward. Alternatively, `.into_data()` is recommended for one-time use. Let's look at a couple of examples for initializing a tensor from different inputs.

```
// Initialization from a given Backend (Wgpu)
let tensor_1 = Tensor::::from_data([1.0, 2.0, 3.0], &device);

// Initialization from a generic Backend
let tensor_2 = Tensor::::from_data(TensorData::from([1.0, 2.0, 3.0]), &device);

// Initialization using from_floats (Recommended for f32 ElementType)
// Will be converted to TensorData internally.
let tensor_3 = Tensor::::from_floats([1.0, 2.0, 3.0], &device);

// Initialization of Int Tensor from array slices
let arr: [i32; 6] = [1, 2, 3, 4, 5, 6];
let tensor_4 = Tensor::::from_data(TensorData::from(&arr[0..3]), &device);

// Initialization from a custom type

struct BodyMetrics {
    age: i8,
    height: i16,
    weight: f32
}

let bmi = BodyMetrics{
    age: 25,
    height: 180,
    weight: 80.0
};
let data = TensorData::from([bmi.age as f32, bmi.height as f32, bmi.weight]);
let tensor_5 = Tensor::::from_data(data, &device);
```

## Ownership and Cloning

Almost all Burn operations take ownership of the input tensors. Therefore, reusing a tensor multiple times will necessitate cloning it. Let's look at an example to understand the ownership rules and cloning better. Suppose we want to do a simple min-max normalization of an input tensor.

```
let input = Tensor::::from_floats([1.0, 2.0, 3.0, 4.0], &device);
let min = input.min();
let max = input.max();
let input = (input - min).div(max - min);
```

With PyTorch tensors, the above code would work as expected. However, Rust's strict

ownership rules will give an error and prevent using the input tensor after the first `.min()` operation. The ownership of the input tensor is transferred to the variable `min` and the input tensor is no longer available for further operations. Burn Tensors like most complex primitives do not implement the `copy` trait and therefore have to be cloned explicitly. Now let's rewrite a working example of doing min-max normalization with cloning.

```
let input = Tensor::<Wgpu, 1>::from_floats([1.0, 2.0, 3.0, 4.0], &device);
let min = input.clone().min();
let max = input.clone().max();
let input = (input.clone() - min.clone()).div(max - min);
println!("{}", input.to_data()); // Success: [0.0, 0.33333334, 0.6666667, 1.0]

// Notice that max, min have been moved in last operation so
// the below print will give an error.
// If we want to use them for further operations,
// they will need to be cloned in similar fashion.
// println!("{}", min.to_data());
```

We don't need to be worried about memory overhead because with cloning, the tensor's buffer isn't copied, and only a reference to it is increased. This makes it possible to determine exactly how many times a tensor is used, which is very convenient for reusing tensor buffers or even fusing operations into a single kernel (*burn-fusion*). For that reason, we don't provide explicit inplace operations. If a tensor is used only one time, inplace operations will always be used when available.

## Tensor Operations

Normally with PyTorch, explicit inplace operations aren't supported during the backward pass, making them useful only for data preprocessing or inference-only model implementations. With Burn, you can focus more on *what* the model should do, rather than on *how* to do it. We take the responsibility of making your code run as fast as possible during training as well as inference. The same principles apply to broadcasting; all operations support broadcasting unless specified otherwise.

Here, we provide a list of all supported operations along with their PyTorch equivalents. Note that for the sake of simplicity, we ignore type signatures. For more details, refer to the [full documentation](#).

### Basic Operations

Those operations are available for all tensor kinds: `Int` , `Float` , and `Bool` .

Burn	PyTorch Equivalent
<code>Tensor::cat(tensors, dim)</code>	<code>torch.cat(tensors, dim)</code>
<code>Tensor::empty(shape, device)</code>	<code>torch.empty(shape, device=device)</code>
<code>Tensor::from_primitive(primitive)</code>	N/A
<code>Tensor::stack(tensors, dim)</code>	<code>torch.stack(tensors, dim)</code>
<code>tensor.all()</code>	<code>tensor.all()</code>
<code>tensor.all_dim(dim)</code>	<code>tensor.all(dim)</code>
<code>tensor.any()</code>	<code>tensor.any()</code>
<code>tensor.any_dim(dim)</code>	<code>tensor.any(dim)</code>
<code>tensor.chunk(num_chunks, dim)</code>	<code>tensor.chunk(num_chunks, dim)</code>
<code>tensor.device()</code>	<code>tensor.device</code>
<code>tensor.dims()</code>	<code>tensor.size()</code>
<code>tensor.equal(other)</code>	<code>x == y</code>
<code>tensor.expand(shape)</code>	<code>tensor.expand(shape)</code>
<code>tensor.flatten(start_dim, end_dim)</code>	<code>tensor.flatten(start_dim, end_dim)</code>
<code>tensor.flip(axes)</code>	<code>tensor.flip(axes)</code>
<code>tensor.into_data()</code>	N/A
<code>tensor.into_primitive()</code>	N/A
<code>tensor.into_scalar()</code>	<code>tensor.item()</code>
<code>tensor.narrow(dim, start, length)</code>	<code>tensor.narrow(dim, start, length)</code>
<code>tensor.not_equal(other)</code>	<code>x != y</code>
<code>tensor.permute(axes)</code>	<code>tensor.permute(axes)</code>
<code>tensor.movedim(src, dst)</code>	<code>tensor.movedim(src, dst)</code>
<code>tensor.repeat_dim(dim, times)</code>	<code>tensor.repeat(*[times if i == dim else 1 for i in range(tensor.dim())])</code>
<code>tensor.repeat(sizes)</code>	<code>tensor.repeat(sizes)</code>
<code>tensor.reshape(shape)</code>	<code>tensor.view(shape)</code>
<code>tensor.shape()</code>	<code>tensor.shape</code>
<code>tensor.slice(ranges)</code>	<code>tensor[(*ranges,)]</code>
<code>tensor.slice_assign(ranges, values)</code>	<code>tensor[(*ranges,)] = values</code>



Burn	PyTorch Equivalent
<code>tensor.squeeze(dim)</code>	<code>tensor.squeeze(dim)</code>
<code>tensor.to_data()</code>	N/A
<code>tensor.to_device(device)</code>	<code>tensor.to(device)</code>
<code>tensor.unsqueeze()</code>	<code>tensor.unsqueeze(0)</code>
<code>tensor.unsqueeze_dim(dim)</code>	<code>tensor.unsqueeze(dim)</code>

## Numeric Operations

Those operations are available for numeric tensor kinds: `Float` and `Int`.

Burn	PyTorch Equivalent
<code>Tensor::eye(size, device)</code>	<code>torch.eye(size, device=device)</code>
<code>Tensor::full(shape, fill_value, device)</code>	<code>torch.full(shape, fill_value, device=device)</code>
<code>Tensor::ones(shape, device)</code>	<code>torch.ones(shape, device=device)</code>
<code>Tensor::zeros(shape)</code>	<code>torch.zeros(shape)</code>
<code>Tensor::zeros(shape, device)</code>	<code>torch.zeros(shape, device=device)</code>
<code>tensor.abs()</code>	<code>torch.abs(tensor)</code>
<code>tensor.add(other)</code> OR <code>tensor + other</code>	<code>tensor + other</code>
<code>tensor.add_scalar(scalar)</code> OR <code>tensor + scalar</code>	<code>tensor + scalar</code>
<code>tensor.all_close(other, atol, rtol)</code>	<code>torch.allclose(tensor, other, atol, rtol)</code>
<code>tensor.argmax(dim)</code>	<code>tensor.argmax(dim)</code>
<code>tensor.argmin(dim)</code>	<code>tensor.argmin(dim)</code>
<code>tensor.argsort(dim)</code>	<code>tensor.argsort(dim)</code>
<code>tensor.argsort_descending(dim)</code>	<code>tensor.argsort(dim, descending=True)</code>
<code>tensor.bool()</code>	<code>tensor.bool()</code>
<code>tensor.clamp(min, max)</code>	<code>torch.clamp(tensor, min=min, max=max)</code>

Burn	PyTorch Equivalent
<code>tensor.clamp_max(max)</code>	<code>torch.clamp(tensor, max=max)</code>
<code>tensor.clamp_min(min)</code>	<code>torch.clamp(tensor, min=min)</code>
<code>tensor.contains_nan()</code>	N/A
<code>tensor.div(other)</code> OR <code>tensor / other</code>	<code>tensor / other</code>
<code>tensor.div_scalar(scalar)</code> OR <code>tensor / scalar</code>	<code>tensor / scalar</code>
<code>tensor.equal_elem(other)</code>	<code>tensor.eq(other)</code>
<code>tensor.full_like(fill_value)</code>	<code>torch.full_like(tensor, fill_value)</code>
<code>tensor.gather(dim, indices)</code>	<code>torch.gather(tensor, dim, indices)</code>
<code>tensor.greater(other)</code>	<code>tensor.gt(other)</code>
<code>tensor.greater_elem(scalar)</code>	<code>tensor.gt(scalar)</code>
<code>tensor.greater_equal(other)</code>	<code>tensor.ge(other)</code>
<code>tensor.greater_equal_elem(scalar)</code>	<code>tensor.ge(scalar)</code>
<code>tensor.is_close(other, atol, rtol)</code>	<code>torch.isclose(tensor, other, atol, rtol)</code>
<code>tensor.is_nan()</code>	<code>torch.isnan(tensor)</code>
<code>tensor.lower(other)</code>	<code>tensor.lt(other)</code>
<code>tensor.lower_elem(scalar)</code>	<code>tensor.lt(scalar)</code>
<code>tensor.lower_equal(other)</code>	<code>tensor.le(other)</code>
<code>tensor.lower_equal_elem(scalar)</code>	<code>tensor.le(scalar)</code>
<code>tensor.mask_fill(mask, value)</code>	<code>tensor.masked_fill(mask, value)</code>
<code>tensor.mask_where(mask, value_tensor)</code>	<code>torch.where(mask, value_tensor, tensor)</code>
<code>tensor.max()</code>	<code>tensor.max()</code>
<code>tensor.max_dim(dim)</code>	<code>tensor.max(dim, keepdim=True)</code>
<code>tensor.max_dim_with_indices(dim)</code>	N/A
<code>tensor.max_pair(other)</code>	<code>torch.Tensor.max(a,b)</code>
<code>tensor.mean()</code>	<code>tensor.mean()</code>
<code>tensor.mean_dim(dim)</code>	<code>tensor.mean(dim, keepdim=True)</code>
<code>tensor.min()</code>	<code>tensor.min()</code>

Burn	PyTorch Equivalent
<code>tensor.min_dim(dim)</code>	<code>tensor.min(dim, keepdim=True)</code>
<code>tensor.min_dim_with_indices(dim)</code>	N/A
<code>tensor.min_pair(other)</code>	<code>torch.Tensor.min(a,b)</code>
<code>tensor.mul(other)</code> OR <code>tensor * other</code>	<code>tensor * other</code>
<code>tensor.mul_scalar(scalar)</code> OR <code>tensor * scalar</code>	<code>tensor * scalar</code>
<code>tensor.neg()</code> OR <code>-tensor</code>	<code>-tensor</code>
<code>tensor.not_equal_elem(scalar)</code>	<code>tensor.ne(scalar)</code>
<code>tensor.ones_like()</code>	<code>torch.ones_like(tensor)</code>
<code>tensor.pad(pads, value)</code>	<code>torch.nn.functional.pad(input, pad, value)</code>
<code>tensor.powf(other)</code> OR <code>tensor.powi(intoother)</code>	<code>tensor.pow(other)</code>
<code>tensor.powf_scalar(scalar)</code> OR <code>tensor.powi_scalar(intscalar)</code>	<code>tensor.pow(scalar)</code>
<code>tensor.prod()</code>	<code>tensor.prod()</code>
<code>tensor.prod_dim(dim)</code>	<code>tensor.prod(dim, keepdim=True)</code>
<code>tensor.rem(other)</code> OR <code>tensor % other</code>	<code>tensor % other</code>
<code>tensor.scatter(dim, indices, values)</code>	<code>tensor.scatter_add(dim, indices, values)</code>
<code>tensor.select(dim, indices)</code>	<code>tensor.index_select(dim, indices)</code>
<code>tensor.select_assign(dim, indices, values)</code>	N/A
<code>tensor.sign()</code>	<code>tensor.sign()</code>
<code>tensor.sort(dim)</code>	<code>tensor.sort(dim).values</code>
<code>tensor.sort_descending(dim)</code>	<code>tensor.sort(dim, descending=True).values</code>
<code>tensor.sort_descending_with_indices(dim)</code>	<code>tensor.sort(dim, descending=True)</code>
<code>tensor.sort_with_indices(dim)</code>	<code>tensor.sort(dim)</code>
<code>tensor.sub(other)</code> OR <code>tensor - other</code>	<code>tensor - other</code>

Burn	PyTorch Equivalent
<code>tensor.sub_scalar(scalar)</code> or <code>tensor - scalar</code>	<code>tensor - scalar</code>
<code>tensor.sum()</code>	<code>tensor.sum()</code>
<code>tensor.sum_dim(dim)</code>	<code>tensor.sum(dim, keepdim=True)</code>
<code>tensor.topk(k, dim)</code>	<code>tensor.topk(k, dim).values</code>
<code>tensor.topk_with_indices(k, dim)</code>	<code>tensor.topk(k, dim)</code>
<code>tensor.tril(diagonal)</code>	<code>torch.tril(tensor, diagonal)</code>
<code>tensor.triu(diagonal)</code>	<code>torch.triu(tensor, diagonal)</code>
<code>tensor.zeros_like()</code>	<code>torch.zeros_like(tensor)</code>

## Float Operations

Those operations are only available for `Float` tensors.

Burn API	PyTorch Equivalent
<code>Tensor::one_hot(index, num_classes, device)</code>	N/A
<code>tensor.ceil()</code>	<code>tensor.ceil()</code>
<code>tensor.cos()</code>	<code>tensor.cos()</code>
<code>tensor.erf()</code>	<code>tensor.erf()</code>
<code>tensor.exp()</code>	<code>tensor.exp()</code>
<code>tensor.floor()</code>	<code>tensor.floor()</code>
<code>tensor.from_floats(floats, device)</code>	N/A
<code>tensor.from_full_precision(tensor)</code>	N/A
<code>tensor.int()</code>	Similar to <code>tensor.to(torch.long)</code>
<code>tensor.log()</code>	<code>tensor.log()</code>
<code>tensor.log1p()</code>	<code>tensor.log1p()</code>
<code>tensor.matmul(other)</code>	<code>tensor.matmul(other)</code>
<code>tensor.random(shape, distribution, device)</code>	N/A
<code>tensor.random_like(distribution)</code>	<code>torch.rand_like()</code> only uniform

Burn API	PyTorch Equivalent
<code>tensor.recip()</code>	<code>tensor.reciprocal()</code>
<code>tensor.round()</code>	<code>tensor.round()</code>
<code>tensor.sin()</code>	<code>tensor.sin()</code>
<code>tensor.sqrt()</code>	<code>tensor.sqrt()</code>
<code>tensor.swap_dims(dim1, dim2)</code>	<code>tensor.transpose(dim1, dim2)</code>
<code>tensor.tanh()</code>	<code>tensor.tanh()</code>
<code>tensor.to_full_precision()</code>	<code>tensor.to(torch.float)</code>
<code>tensor.transpose()</code>	<code>tensor.T</code>
<code>tensor.var(dim)</code>	<code>tensor.var(dim)</code>
<code>tensor.var_bias(dim)</code>	N/A
<code>tensor.var_mean(dim)</code>	N/A
<code>tensor.var_mean_bias(dim)</code>	N/A

## Int Operations

Those operations are only available for `Int` tensors.

Burn API	PyTorch Equivalent
<code>tensor.arange(5..10, device)</code>	<code>tensor.arange(start=5, end=10, device=device)</code>
<code>tensor.arange_step(5..10, 2, device)</code>	<code>tensor.arange(start=5, end=10, step=2, device=device)</code>
<code>tensor.float()</code>	<code>tensor.to(torch.float)</code>
<code>tensor.from_ints(ints)</code>	N/A
<code>tensor.int_random(shape, distribution, device)</code>	N/A
<code>tensor.cartesian_grid(shape, device)</code>	N/A
<code>tensor.one_hot(num_classes)</code>	N/A

## Bool Operations

Those operations are only available for `Bool` tensors.

Burn API	PyTorch Equivalent
<code>Tensor::diag_mask(shape, diagonal)</code>	N/A
<code>Tensor::tril_mask(shape, diagonal)</code>	N/A
<code>Tensor::triu_mask(shape, diagonal)</code>	N/A
<code>tensor.argwhere()</code>	<code>tensor.argwhere()</code>
<code>tensor.float()</code>	<code>tensor.to(torch.float)</code>
<code>tensor.int()</code>	<code>tensor.to(torch.long)</code>
<code>tensor.nonzero()</code>	<code>tensor.nonzero(as_tuple=True)</code>
<code>tensor.not()</code>	<code>tensor.logical_not()</code>

## Quantization Operations

Those operations are only available for `Float` tensors on backends that implement quantization strategies.

Burn API	PyTorch Equivalent
<code>tensor.quantize(scheme, qparams)</code>	N/A
<code>tensor.dequantize()</code>	N/A

## Activation Functions

Burn API	PyTorch Equivalent
<code>activation::gelu(tensor)</code>	<code>nn.functional.gelu(tensor)</code>
<code>activation::hard_sigmoid(tensor, alpha, beta)</code>	<code>nn.functional.hardsigmoid(tensor)</code>
<code>activation::leaky_relu(tensor, negative_slope)</code>	<code>nn.functional.leaky_relu(tensor, negative_slope)</code>
<code>activation::log_sigmoid(tensor)</code>	<code>nn.functional.log_sigmoid(tensor)</code>
<code>activation::log_softmax(tensor, dim)</code>	<code>nn.functional.log_softmax(tensor, dim)</code>
<code>activation::mish(tensor)</code>	<code>nn.functional.mish(tensor)</code>
<code>activation::prelu(tensor, alpha)</code>	<code>nn.functional.prelu(tensor, weight)</code>

Burn API	PyTorch Equivalent
<code>activation::quiet_softmax(tensor, dim)</code>	<code>nn.functional.quiet_softmax(tensor, dim)</code>
<code>activation::relu(tensor)</code>	<code>nn.functional.relu(tensor)</code>
<code>activation::sigmoid(tensor)</code>	<code>nn.functional.sigmoid(tensor)</code>
<code>activation::silu(tensor)</code>	<code>nn.functional.silu(tensor)</code>
<code>activation::softmax(tensor, dim)</code>	<code>nn.functional.softmax(tensor, dim)</code>
<code>activation::softmin(tensor, dim)</code>	<code>nn.functional.softmin(tensor, dim)</code>
<code>activation::softplus(tensor, beta)</code>	<code>nn.functional.softplus(tensor, beta)</code>
<code>activation::tanh(tensor)</code>	<code>nn.functional.tanh(tensor)</code>

## Displaying Tensor Details

Burn provides flexible options for displaying tensor information, allowing you to control the level of detail and formatting to suit your needs.

### Basic Display

To display a detailed view of a tensor, you can simply use Rust's `println!` or `format!` macros:

```
let tensor = Tensor::<Backend, 2>::full([2, 3], 0.123456789,
    &Default::default());
println!("{}", tensor);
```

This will output:

```
Tensor {
  data:
  [[0.12345679, 0.12345679, 0.12345679],
  [0.12345679, 0.12345679, 0.12345679]],
  shape: [2, 3],
  device: Cpu,
  backend: "ndarray",
  kind: "Float",
  dtype: "f32",
}
```

## Controlling Precision

You can control the number of decimal places displayed using Rust's formatting syntax:

```
println!("{:.2}", tensor);
```

Output:

```
Tensor {  
  data:  
  [[0.12, 0.12, 0.12],  
   [0.12, 0.12, 0.12]],  
  shape: [2, 3],  
  device: Cpu,  
  backend: "ndarray",  
  kind: "Float",  
  dtype: "f32",  
}
```

## Global Print Options

For more fine-grained control over tensor printing, Burn provides a `PrintOptions` struct and a `set_print_options` function:

```
use burn::tensor::{set_print_options, PrintOptions};  
  
let print_options = PrintOptions {  
  precision: Some(2),  
  ..Default::default()  
};  
  
set_print_options(print_options);
```

Options:

- `precision` : Number of decimal places for floating-point numbers (default: None)
- `threshold` : Maximum number of elements to display before summarizing (default: 1000)
- `edge_items` : Number of items to show at the beginning and end of each dimension when summarizing (default: 3)



## Checking Tensor Closeness

Burn provides a utility function `check_closeness` to compare two tensors and assess their similarity. This function is particularly useful for debugging and validating tensor operations, especially when working with floating-point arithmetic where small numerical differences can accumulate. It's also valuable when comparing model outputs during the process of importing models from other frameworks, helping to ensure that the imported model produces results consistent with the original.

Here's an example of how to use `check_closeness`:

```
use burn::tensor::{check_closeness, Tensor};
type B = burn::backend::NdArray;

let device = Default::default();
let tensor1 = Tensor::<B, 1>::from_floats(
    [1.0, 2.0, 3.0, 4.0, 5.0, 6.001, 7.002, 8.003, 9.004, 10.1],
    &device,
);
let tensor2 = Tensor::<B, 1>::from_floats(
    [1.0, 2.0, 3.0, 4.000, 5.0, 6.0, 7.001, 8.002, 9.003, 10.004],
    &device,
);

check_closeness(&tensor1, &tensor2);
```

The `check_closeness` function compares the two input tensors element-wise, checking their absolute differences against a range of epsilon values. It then prints a detailed report showing the percentage of elements that are within each tolerance level.

The output provides a breakdown for different epsilon values, allowing you to assess the closeness of the tensors at various precision levels. This is particularly helpful when dealing with operations that may introduce small numerical discrepancies.

The function uses color-coded output to highlight the results:

- **Green [PASS]**: All elements are within the specified tolerance.
- **Yellow [WARN]**: Most elements (90% or more) are within tolerance.
- **Red [FAIL]**: Significant differences are detected.

This utility can be invaluable when implementing or debugging tensor operations,

especially those involving complex mathematical computations or when porting algorithms from other frameworks. It's also an essential tool when verifying the accuracy of imported models, ensuring that the Burn implementation produces results that closely match those of the original model.

# Autodiff

Burn's tensor also supports auto-differentiation, which is an essential part of any deep learning framework. We introduced the `Backend` trait in the [previous section](#), but Burn also has another trait for autodiff: `AutodiffBackend`.

However, not all tensors support auto-differentiation; you need a backend that implements both the `Backend` and `AutodiffBackend` traits. Fortunately, you can add auto-differentiation capabilities to any backend using a backend decorator: `type MyAutodiffBackend = Autodiff<MyBackend>`. This decorator implements both the `AutodiffBackend` and `Backend` traits by maintaining a dynamic computational graph and utilizing the inner backend to execute tensor operations.

The `AutodiffBackend` trait adds new operations on float tensors that can't be called otherwise. It also provides a new associated type, `B::Gradients`, where each calculated gradient resides.

```
fn calculate_gradients<B: AutodiffBackend>(tensor: Tensor<B, 2>) ->
B::Gradients {
    let mut gradients = tensor.clone().backward();

    let tensor_grad = tensor.grad(&gradients);          // get
    let tensor_grad = tensor.grad_remove(&mut gradients); // pop

    gradients
}
```

Note that some functions will always be available even if the backend doesn't implement the `AutodiffBackend` trait. In such cases, those functions will do nothing.

Burn API	PyTorch Equivalent
<code>tensor.detach()</code>	<code>tensor.detach()</code>
<code>tensor.require_grad()</code>	<code>tensor.requires_grad()</code>
<code>tensor.is_require_grad()</code>	<code>tensor.requires_grad</code>
<code>tensor.set_require_grad(require_grad)</code>	<code>tensor.requires_grad(False)</code>

However, you're unlikely to make any mistakes since you can't call `backward` on a tensor that is on a backend that doesn't implement `AutodiffBackend`. Additionally, you can't retrieve the gradient of a tensor without an autodiff backend.

## Difference with PyTorch

The way Burn handles gradients is different from PyTorch. First, when calling `backward`, each parameter doesn't have its `grad` field updated. Instead, the backward pass returns all the calculated gradients in a container. This approach offers numerous benefits, such as the ability to easily send gradients to other threads.

You can also retrieve the gradient for a specific parameter using the `grad` method on a tensor. Since this method takes the gradients as input, it's hard to forget to call `backward` beforehand. Note that sometimes, using `grad_remove` can improve performance by allowing inplace operations.

In PyTorch, when you don't need gradients for inference or validation, you typically need to scope your code using a block.

```
# Inference mode
torch.inference():
  # your code
  ...

# Or no grad
torch.no_grad():
  # your code
  ...
```

With Burn, you don't need to wrap the backend with the `Autodiff` for inference, and you can call `inner()` to obtain the inner tensor, which is useful for validation.

```
/// Use `B: AutodiffBackend`
fn example_validation<B: AutodiffBackend>(tensor: Tensor<B, 2>) {
  let inner_tensor: Tensor<B::InnerBackend, 2> = tensor.inner();
  let _ = inner_tensor + 5;
}

/// Use `B: Backend`
fn example_inference<B: Backend>(tensor: Tensor<B, 2>) {
  let _ = tensor + 5;
  ...
}
```

## Gradients with Optimizers

We've seen how gradients can be used with tensors, but the process is a bit different when working with optimizers from `burn-core`. To work with the `Module` trait, a translation step is required to link tensor parameters with their gradients. This step is necessary to easily

support gradient accumulation and training on multiple devices, where each module can be forked and run on different devices in parallel. We'll explore deeper into this topic in the [Module](#) section.

# Module

The `Module` derive allows you to create your own neural network modules, similar to PyTorch. The `derive` function only generates the necessary methods to essentially act as a parameter container for your type, it makes no assumptions about how the forward pass is declared.

```
use burn::module::Module;
use burn::tensor::backend::Backend;

#[derive(Module, Debug)]
pub struct PositionWiseFeedForward<B: Backend> {
    linear_inner: Linear<B>,
    linear_outer: Linear<B>,
    dropout: Dropout,
    gelu: Gelu,
}

impl<B: Backend> PositionWiseFeedForward<B> {
    /// Normal method added to a struct.
    pub fn forward<const D: usize>(&self, input: Tensor<B, D>) -> Tensor<B, D>
    {
        let x = self.linear_inner.forward(input);
        let x = self.gelu.forward(x);
        let x = self.dropout.forward(x);

        self.linear_outer.forward(x)
    }
}
```

Note that all fields declared in the struct must also implement the `Module` trait.

## Tensor

If you want to create your own module that contains tensors, and not just other modules defined with the `Module` derive, you need to be careful to achieve the behavior you want.

- `Param<Tensor<B, D>>` : If you want the tensor to be included as a parameter of your modules, you need to wrap the tensor in a `Param` struct. This will create an ID that will be used to identify this parameter. This is essential when performing module optimization and when saving states such as optimizer and module checkpoints. Note that a module's record only contains parameters.

- `Param<Tensor<B, D>>.set_require_grad(false)` : If you want the tensor to be included as a parameter of your modules, and therefore saved with the module's weights, but you don't want it to be updated by the optimizer.
- `Tensor<B, D>` : If you want the tensor to act as a constant that can be recreated when instantiating a module. This can be useful when generating sinusoidal embeddings, for example.

## Methods

These methods are available for all modules.

Burn API	PyTorch Equivalent
<code>module.devices()</code>	N/A
<code>module.fork(device)</code>	Similar to <code>module.to(device).detach()</code>
<code>module.to_device(device)</code>	<code>module.to(device)</code>
<code>module.no_grad()</code>	<code>module.require_grad_(False)</code>
<code>module.num_params()</code>	N/A
<code>module.visit(visitor)</code>	N/A
<code>module.map(mapper)</code>	N/A
<code>module.into_record()</code>	Similar to <code>state_dict</code>
<code>module.load_record(record)</code>	Similar to <code>load_state_dict(state_dict)</code>
<code>module.save_file(file_path, recorder)</code>	N/A
<code>module.load_file(file_path, recorder)</code>	N/A

Similar to the backend trait, there is also the `AutodiffModule` trait to signify a module with autodiff support.

Burn API	PyTorch Equivalent
<code>module.valid()</code>	<code>module.eval()</code>

## Visitor & Mapper

As mentioned earlier, modules primarily function as parameter containers. Therefore, we naturally offer several ways to perform functions on each parameter. This is distinct from PyTorch, where extending module functionalities is not as straightforward.

The `map` and `visitor` methods are quite similar but serve different purposes. Mapping is used for potentially mutable operations where each parameter of a module can be updated to a new value. In Burn, optimizers are essentially just sophisticated module mappers. Visitors, on the other hand, are used when you don't intend to modify the module but need to retrieve specific information from it, such as the number of parameters or a list of devices in use.

You can implement your own mapper or visitor by implementing these simple traits:

```
/// Module visitor trait.
pub trait ModuleVisitor<B: Backend> {
    /// Visit a tensor in the module.
    fn visit<const D: usize>(&mut self, id: ParamId, tensor: &Tensor<B, D>);
}

/// Module mapper trait.
pub trait ModuleMapper<B: Backend> {
    /// Map a tensor in the module.
    fn map<const D: usize>(&mut self, id: ParamId, tensor: Tensor<B, D>) ->
        Tensor<B, D>;
}
```

## Module Display

Burn provides a simple way to display the structure of a module and its configuration at a glance. You can print the module to see its structure, which is useful for debugging and tracking changes across different versions of a module. (See the print output of the [Basic Workflow Model](#) example.)

To customize the display of a module, you can implement the `ModuleDisplay` trait for your module. This will change the default display settings for the module and its children. Note that `ModuleDisplay` is automatically implemented for all modules, but you can override it to customize the display by annotating the module with `#[module(custom_display)]`.



```

#[derive(Module, Debug)]
#[module(custom_display)]
pub struct PositionWiseFeedForward<B: Backend> {
    linear_inner: Linear<B>,
    linear_outer: Linear<B>,
    dropout: Dropout,
    gelu: Gelu,
}

impl<B: Backend> ModuleDisplay for PositionWiseFeedForward<B> {
    /// Custom settings for the display of the module.
    /// If `None` is returned, the default settings will be used.
    fn custom_settings(&self) -> Option<burn::module::DisplaySettings> {
        DisplaySettings::new()
            // Will show all attributes (default is false)
            .with_show_all_attributes(false)
            // Will show each attribute on a new line (default is true)
            .with_new_line_after_attribute(true)
            // Will show the number of parameters (default is true)
            .with_show_num_parameters(true)
            // Will indent by 2 spaces (default is 2)
            .with_indentation_size(2)
            // Will show the parameter ID (default is false)
            .with_show_param_id(false)
            // Convenience method to wrap settings in Some()
            .optional()
    }

    /// Custom content to be displayed.
    /// If `None` is returned, the default content will be used
    /// (all attributes of the module)
    fn custom_content(&self, content: Content) -> Option<Content> {
        content
            .add("linear_inner", &self.linear_inner)
            .add("linear_outer", &self.linear_outer)
            .add("anything", "anything_else")
            .optional()
    }
}

```

## Built-in Modules

Burn comes with built-in modules that you can use to build your own modules.

### General

Burn API	PyTorch Equivalent
BatchNorm	<code>nn.BatchNorm1d</code> , <code>nn.BatchNorm2d</code> etc.
Dropout	<code>nn.Dropout</code>
Embedding	<code>nn.Embedding</code>
Gelu	<code>nn.Gelu</code>
GroupNorm	<code>nn.GroupNorm</code>
HardSigmoid	<code>nn.Hardsigmoid</code>
InstanceNorm	<code>nn.InstanceNorm1d</code> , <code>nn.InstanceNorm2d</code> etc.
LayerNorm	<code>nn.LayerNorm</code>
LeakyRelu	<code>nn.LeakyReLU</code>
Linear	<code>nn.Linear</code>
Prelu	<code>nn.PReLU</code>
Relu	<code>nn.ReLU</code>
RmsNorm	<i>No direct equivalent</i>
SwiGlu	<i>No direct equivalent</i>
Interpolate1d	<i>No direct equivalent</i>
Interpolate2d	<i>No direct equivalent</i>

## Convolutions

Burn API	PyTorch Equivalent
Conv1d	<code>nn.Conv1d</code>
Conv2d	<code>nn.Conv2d</code>
Conv3d	<code>nn.Conv3d</code>
ConvTranspose1d	<code>nn.ConvTranspose1d</code>
ConvTranspose2d	<code>nn.ConvTranspose2d</code>
ConvTranspose3d	<code>nn.ConvTranspose3d</code>
DeformConv2d	<code>torchvision.ops.DeformConv2d</code>

## Pooling

Burn API	PyTorch Equivalent
----------	--------------------

Burn API	PyTorch Equivalent
AdaptiveAvgPool1d	<code>nn.AdaptiveAvgPool1d</code>
AdaptiveAvgPool2d	<code>nn.AdaptiveAvgPool2d</code>
AvgPool1d	<code>nn.AvgPool1d</code>
AvgPool2d	<code>nn.AvgPool2d</code>
MaxPool1d	<code>nn.MaxPool1d</code>
MaxPool2d	<code>nn.MaxPool2d</code>

## RNNs

Burn API	PyTorch Equivalent
Gru	<code>nn.GRU</code>
Lstm / BiLstm	<code>nn.LSTM</code>
GateController	<i>No direct equivalent</i>

## Transformer

Burn API	PyTorch Equivalent
MultiHeadAttention	<code>nn.MultiheadAttention</code>
TransformerDecoder	<code>nn.TransformerDecoder</code>
TransformerEncoder	<code>nn.TransformerEncoder</code>
PositionalEncoding	<i>No direct equivalent</i>
RotaryEncoding	<i>No direct equivalent</i>

## Loss

Burn API	PyTorch Equivalent
CrossEntropyLoss	<code>nn.CrossEntropyLoss</code>
MseLoss	<code>nn.MSELoss</code>
HuberLoss	<code>nn.HuberLoss</code>

# Learner

The [burn-train](#) crate encapsulates multiple utilities for training deep learning models. The goal of the crate is to provide users with a well-crafted and flexible training loop, so that projects do not have to write such components from the ground up. Most of the interactions with `burn-train` will be with the `LearnerBuilder` struct, briefly presented in the previous [training section](#). This struct enables you to configure the training loop, offering support for registering metrics, enabling logging, checkpointing states, using multiple devices, and so on.

There are still some assumptions in the current provided APIs, which may make them inappropriate for your learning requirements. Indeed, they assume your model will learn from a training dataset and be validated against another dataset. This is the most common paradigm, allowing users to do both supervised and unsupervised learning as well as fine-tuning. However, for more complex requirements, creating a [custom training loop](#) might be what you need.

## Usage

The learner builder provides numerous options when it comes to configurations.

Configuration	Description
Training Metric	Register a training metric
Validation Metric	Register a validation metric
Training Metric Plot	Register a training metric with plotting (requires the metric to be numeric)
Validation Metric Plot	Register a validation metric with plotting (requires the metric to be numeric)
Metric Logger	Configure the metric loggers (default is saving them to files)
Renderer	Configure how to render metrics (default is CLI)
Grad Accumulation	Configure the number of steps before applying gradients
File Checkpointer	Configure how the model, optimizer and scheduler states are saved
Num Epochs	Set the number of epochs.
Devices	Set the devices to be used
Checkpoint	Restart training from a checkpoint

Configuration	Description
Application logging	Configure the application logging installer (default is writing to <code>experiment.log</code> )

When the builder is configured at your liking, you can then move forward to build the learner. The build method requires three inputs: the model, the optimizer and the learning rate scheduler. Note that the latter can be a simple float if you want it to be constant during training.

The result will be a newly created `Learner` struct, which has only one method, the `fit` function which must be called with the training and validation dataloaders. This will start the training and return the trained model once finished.

Again, please refer to the [training section](#) for a relevant code snippet.

## Artifacts

When creating a new builder, all the collected data will be saved under the directory provided as the argument to the `new` method. Here is an example of the data layout for a model recorded using the compressed message pack format, with the accuracy and loss metrics registered:

```
├── experiment.log
├── checkpoint
│   ├── model-1.mpk.gz
│   ├── optim-1.mpk.gz
│   ├── scheduler-1.mpk.gz
│   ├── model-2.mpk.gz
│   ├── optim-2.mpk.gz
│   └── scheduler-2.mpk.gz
├── train
│   ├── epoch-1
│   │   ├── Accuracy.log
│   │   └── Loss.log
│   └── epoch-2
│       ├── Accuracy.log
│       └── Loss.log
└── valid
    ├── epoch-1
    │   ├── Accuracy.log
    │   └── Loss.log
    └── epoch-2
        ├── Accuracy.log
        └── Loss.log
```

You can choose to save or synchronize that local directory with a remote file system, if desired. The file checkpointer is capable of automatically deleting old checkpoints according to a specified configuration.

# Metric

When working with the learner, you have the option to record metrics that will be monitored throughout the training process. We currently offer a restricted range of metrics.

Metric	Description
Accuracy	Calculate the accuracy in percentage
TopKAccuracy	Calculate the top-k accuracy in percentage
Loss	Output the loss used for the backward pass
CPU Temperature	Fetch the temperature of CPUs
CPU Usage	Fetch the CPU utilization
CPU Memory Usage	Fetch the CPU RAM usage
GPU Temperature	Fetch the GPU temperature
Learning Rate	Fetch the current learning rate for each optimizer step
CUDA	Fetch general CUDA metrics such as utilization

In order to use a metric, the output of your training step has to implement the `Adaptor` trait from `burn-train::metric`. Here is an example for the classification output, already provided with the crate.

```

/// Simple classification output adapted for multiple metrics.
#[derive(new)]
pub struct ClassificationOutput<B: Backend> {
    /// The loss.
    pub loss: Tensor<B, 1>,

    /// The output.
    pub output: Tensor<B, 2>,

    /// The targets.
    pub targets: Tensor<B, 1, Int>,
}

impl<B: Backend> Adaptor<AccuracyInput<B>> for ClassificationOutput<B> {
    fn adapt(&self) -> AccuracyInput<B> {
        AccuracyInput::new(self.output.clone(), self.targets.clone())
    }
}

impl<B: Backend> Adaptor<LossInput<B>> for ClassificationOutput<B> {
    fn adapt(&self) -> LossInput<B> {
        LossInput::new(self.loss.clone())
    }
}

```

## Custom Metric

Generating your own custom metrics is done by implementing the `Metric` trait.

```

/// Metric trait.
///
/// Implementations should define their own input type only used by the metric.
/// This is important since some conflict may happen when the model output is
/// adapted for each
/// metric's input type.
///
/// The only exception is for metrics that don't need any input, setting the
/// associated type
/// to the null type `()`.
pub trait Metric: Send + Sync {
    /// The input type of the metric.
    type Input;

    /// Updates the metric state and returns the current metric entry.
    fn update(&mut self, item: &Self::Input, metadata: &MetricMetadata) ->
MetricEntry;
    /// Clear the metric state.
    fn clear(&mut self);
}

```



As an example, let's see how the loss metric is implemented.

```
/// The loss metric.
#[derive(Default)]
pub struct LossMetric<B: Backend> {
    state: NumericMetricState,
    _b: B,
}

/// The loss metric input type.
#[derive(new)]
pub struct LossInput<B: Backend> {
    tensor: Tensor<B, 1>,
}

impl<B: Backend> Metric for LossMetric<B> {
    type Input = LossInput<B>;

    fn update(&mut self, loss: &Self::Input, _metadata: &MetricMetadata) ->
MetricEntry {
        let loss = loss.tensor.clone().mean().into_scalar().elem::<f64>();

        self.state
            .update(loss, 1, FormatOptions::new("Loss").precision(2))
    }

    fn clear(&mut self) {
        self.state.reset()
    }
}
```

When the metric you are implementing is numeric in nature, you may want to also implement the `Numeric` trait. This will allow your metric to be plotted.

```
impl<B: Backend> Numeric for LossMetric<B> {
    fn value(&self) -> f64 {
        self.state.value()
    }
}
```

# Config

When writing scientific code, you normally have a lot of values that are set, and Deep Learning is no exception. Python has the possibility to define default parameters for functions, which helps improve the developer experience. However, this has the downside of potentially breaking your code when upgrading to a new version, as the default values might change without your knowledge, making debugging very challenging.

With that in mind, we came up with the Config system. It's a simple Rust derive that you can apply to your types, allowing you to define default values with ease. Additionally, all configs can be serialized, reducing potential bugs when upgrading versions and improving reproducibility.

```
use burn::config::Config;

#[derive(Config)]
pub struct MyModuleConfig {
    d_model: usize,
    d_ff: usize,
    #[config(default = 0.1)]
    dropout: f64,
}
```

The derive also adds useful `with_` methods for every attribute of your config, similar to a builder pattern, along with a `save` method.

```
fn main() {
    let config = MyModuleConfig::new(512, 2048);
    println!("{}", config.d_model); // 512
    println!("{}", config.d_ff); // 2048
    println!("{}", config.dropout); // 0.1
    let config = config.with_dropout(0.2);
    println!("{}", config.dropout); // 0.2

    config.save("config.json").unwrap();
}
```

## Good practices

By using the config type it is easy to create new module instances. The initialization method should be implemented on the config type with the device as argument.

```
impl MyModuleConfig {  
    /// Create a module on the given device.  
    pub fn init<B: Backend>(&self, device: &B::Device) -> MyModule {  
        MyModule {  
            linear: LinearConfig::new(self.d_model, self.d_ff).init(device),  
            dropout: DropoutConfig::new(self.dropout).init(),  
        }  
    }  
}
```

Then we could add this line to the above `main`:

```
use burn::backend::Wgpu;  
let device = Default::default();  
let my_module = config.init::<Wgpu>(&device);
```

# Record

Records are how states are saved with Burn. Compared to most other frameworks, Burn has its own advanced saving mechanism that allows interoperability between backends with minimal possible runtime errors. There are multiple reasons why Burn decided to create its own saving formats.

First, Rust has [serde](#), which is an extremely well-developed serialization and deserialization library that also powers the `safetensors` format developed by Hugging Face. If used properly, all the validations are done when deserializing, which removes the need to write validation code. Since modules in Burn are created with configurations, they can't implement serialization and deserialization. That's why the record system was created: allowing you to save the state of modules independently of the backend in use extremely fast while still giving you all the flexibility possible to include any non-serializable field within your module.

## Why not use `safetensors`?

[safetensors](#) uses `serde` with the JSON file format and only supports serializing and deserializing tensors. The record system in Burn gives you the possibility to serialize any type, which is very useful for optimizers that save their state, but also for any non-standard, cutting-edge modeling needs you may have. Additionally, the record system performs automatic precision conversion by using Rust types, making it more reliable with fewer manual manipulations.

It is important to note that the `safetensors` format uses the word *safe* to distinguish itself from Pickle, which is vulnerable to Python code injection. On our end, the simple fact that we use Rust already ensures that no code injection is possible. If your storage mechanism doesn't handle data corruption, you might prefer a recorder that performs checksum validation (i.e., any recorder with Gzip compression).

## Recorder

Recorders are independent of the backend and serialize records with precision and a format. Note that the format can also be in-memory, allowing you to save the records directly into bytes.

Recorder	Format	Compression
DefaultFileRecorder	File - Named MessagePack	None

Recorder	Format	Compression
NamedMpkFileRecorder	File - Named MessagePack	None
NamedMpkGzFileRecorder	File - Named MessagePack	Gzip
BinFileRecorder	File - Binary	None
BinGzFileRecorder	File - Binary	Gzip
JsonGzFileRecorder	File - Json	Gzip
PrettyJsonFileRecorder	File - Pretty Json	Gzip
BinBytesRecorder	In Memory - Binary	None

Each recorder supports precision settings decoupled from the precision used for training or inference. These settings allow you to define the floating-point and integer types that will be used for serialization and deserialization.

Setting	Float Precision	Integer Precision
DoublePrecisionSettings	f64	i64
FullPrecisionSettings	f32	i32
HalfPrecisionSettings	f16	i16

Note that when loading a record into a module, the type conversion is automatically handled, so you can't encounter errors. The only crucial aspect is using the same recorder for both serialization and deserialization; otherwise, you will encounter loading errors.

### Which recorder should you use?

- If you want fast serialization and deserialization, choose a recorder without compression. The one with the lowest file size without compression is the binary format; otherwise, the named MessagePack could be used.
- If you want to save models for storage, you can use compression, but avoid using the binary format, as it may not be backward compatible.
- If you want to debug your model's weights, you can use the pretty JSON format.
- If you want to deploy with `no-std`, use the in-memory binary format and include the bytes with the compiled code.

For examples on saving and loading records, take a look at [Saving and Loading Models](#).

# Dataset

At its core, a dataset is a collection of data typically related to a specific analysis or processing task. The data modality can vary depending on the task, but most datasets primarily consist of images, texts, audio or videos.

This data source represents an integral part of machine learning to successfully train a model. Thus, it is essential to provide a convenient and performant API to handle your data. Since this process varies wildly from one problem to another, it is defined as a trait that should be implemented on your type. The dataset trait is quite similar to the dataset abstract class in PyTorch:

```
pub trait Dataset<I>: Send + Sync {  
    fn get(&self, index: usize) -> Option<I>;  
    fn len(&self) -> usize;  
}
```

The dataset trait assumes a fixed-length set of items that can be randomly accessed in constant time. This is a major difference from datasets that use Apache Arrow underneath to improve streaming performance. Datasets in Burn don't assume *how* they are going to be accessed; it's just a collection of items.

However, you can compose multiple dataset transformations to lazily obtain what you want with zero pre-processing, so that your training can start instantly!

## Transformation

Transformations in Burn are all lazy and modify one or multiple input datasets. The goal of these transformations is to provide you with the necessary tools so that you can model complex data distributions.

Transformation	Description
SamplerDataset	Samples items from a dataset. This is a convenient way to model a dataset as a probability distribution of a fixed size.
ShuffledDataset	Maps each input index to a random index, similar to a dataset sampled without replacement.
PartialDataset	Returns a view of the input dataset with a specified range.
MapperDataset	Computes a transformation lazily on the input dataset.

Transformation	Description
ComposedDataset	Composes multiple datasets together to create a larger one without copying any data.
WindowsDataset	Dataset designed to work with overlapping windows of data extracted from an input dataset.

Let us look at the basic usages of each dataset transform and how they can be composed together. These transforms are lazy by default except when specified, reducing the need for unnecessary intermediate allocations and improving performance. The full documentation of each transform can be found at the [API reference](#).

- **SamplerDataset:** This transform can be used to sample items from a dataset with (default) or without replacement. Transform is initialized with a sampling size which can be bigger or smaller than the input dataset size. This is particularly useful in cases where we want to checkpoint larger datasets more often during training and smaller datasets less often as the size of an epoch is now controlled by the sampling size. Sample usage:

```
type DbPedia = SqliteDataset<DbPediaItem>;  
let dataset: DbPedia = HuggingfaceDatasetLoader::new("dbpedia_14")  
    .dataset("train").  
    .unwrap();  
  
let dataset = SamplerDataset<DbPedia, DbPediaItem>::new(dataset, 10000);
```

- **ShuffledDataset:** This transform can be used to shuffle the items of a dataset. Particularly useful before splitting the raw dataset into train/test splits. Can be initialized with a seed to ensure reproducibility.

```
let dataset = ShuffledDataset<DbPedia, DbPediaItem>::with_seed(dataset, 42);
```

- **PartialDataset:** This transform is useful to return a view of the dataset with specified start and end indices. Used to create train/val/test splits. In the example below, we show how to chain ShuffledDataset and PartialDataset to create splits.

```
// define chained dataset type here for brevity
type PartialData = PartialDataset<ShuffledDataset<DbPedia, DbPediaItem>>;
let len = dataset.len();
let split == "train"; // or "val"/"test"

let data_split = match split {
    "train" => PartialData::new(dataset, 0, len * 8 / 10), // Get first
80% dataset
    "test" => PartialData::new(dataset, len * 8 / 10, len), // Take
remaining 20%
    _ => panic!("Invalid split type"), // Handle
unexpected split types
};
```

- **MapperDataset:** This transform is useful to apply a transformation on each of the items of a dataset. Particularly useful for normalization of image data when channel means are known.
- **ComposedDataset:** This transform is useful to compose multiple datasets downloaded from multiple sources (say different HuggingfaceDatasetLoader sources) into a single bigger dataset which can be sampled from one source.
- **WindowsDataset:** This transform is useful to create overlapping windows of a dataset. Particularly useful for sequential Time series Data, for example when working with an LSTM.

## Storage

There are multiple dataset storage options available for you to choose from. The choice of the dataset to use should be based on the dataset's size as well as its intended purpose.

Storage	Description
InMemDataset	In-memory dataset that uses a vector to store items. Well-suited for smaller datasets.
SqliteDataset	Dataset that uses <a href="#">SQLite</a> to index items that can be saved in a simple SQL database file. Well-suited for larger datasets.
DataframeDataset	Dataset that uses <a href="#">Polars</a> dataframe to store and manage data. Well-suited for efficient data manipulation and analysis.



## Sources

For now, there are only a couple of dataset sources available with Burn, but more to come!

### Hugging Face

You can easily import any Hugging Face dataset with Burn. We use SQLite as the storage to avoid downloading the model each time or starting a Python process. You need to know the format of each item in the dataset beforehand. Here's an example with the [dbpedia dataset](#).

```
#[derive(Clone, Debug, serde::Serialize, serde::Deserialize)]
pub struct DbPediaItem {
    pub title: String,
    pub content: String,
    pub label: usize,
}

fn main() {
    let dataset: SqliteDataset<DbPediaItem> =
    HuggingfaceDatasetLoader::new("dbpedia_14")
        .dataset("train") // The training split.
        .unwrap();
}
```

We see that items must derive `serde::Serialize`, `serde::Deserialize`, `Clone`, and `Debug`, but those are the only requirements.

### Images

`ImageFolderDataset` is a generic vision dataset used to load images from disk. It is currently available for multi-class and multi-label classification tasks.

```
// Create an image classification dataset from the root folder,
// where images for each class are stored in their respective folder.
//
// For example:
// root/dog/dog1.png
// root/dog/dog2.png
// ...
// root/cat/cat1.png
let dataset = ImageFolderDataset::new_classification("path/to/dataset/
root").unwrap();
```

```
// Create a multi-label image classification dataset from a list of items,
// where each item is a tuple `(image path, labels)`, and a list of classes
// in the dataset.
//
// For example:
let items = vec![
  ("root/dog/dog1.png", vec!["animal".to_string(), "dog".to_string()]),
  ("root/cat/cat1.png", vec!["animal".to_string(), "cat".to_string()]),
];
let dataset = ImageFolderDataset::new_multilabel_classification_with_items(
  items,
  &["animal", "cat", "dog"],
)
.unwrap();
```

## Comma-Separated Values (CSV)

Loading records from a simple CSV file in-memory is simple with the `InMemDataset`:

```
// Build dataset from csv with tab ('\t') delimiter.
// The reader can be configured for your particular file.
let mut rdr = csv::ReaderBuilder::new();
let rdr = rdr.delimiter(b'\t');

let dataset = InMemDataset::from_csv("path/to/csv", rdr).unwrap();
```

Note that this requires the `csv` crate.

### What about streaming datasets?

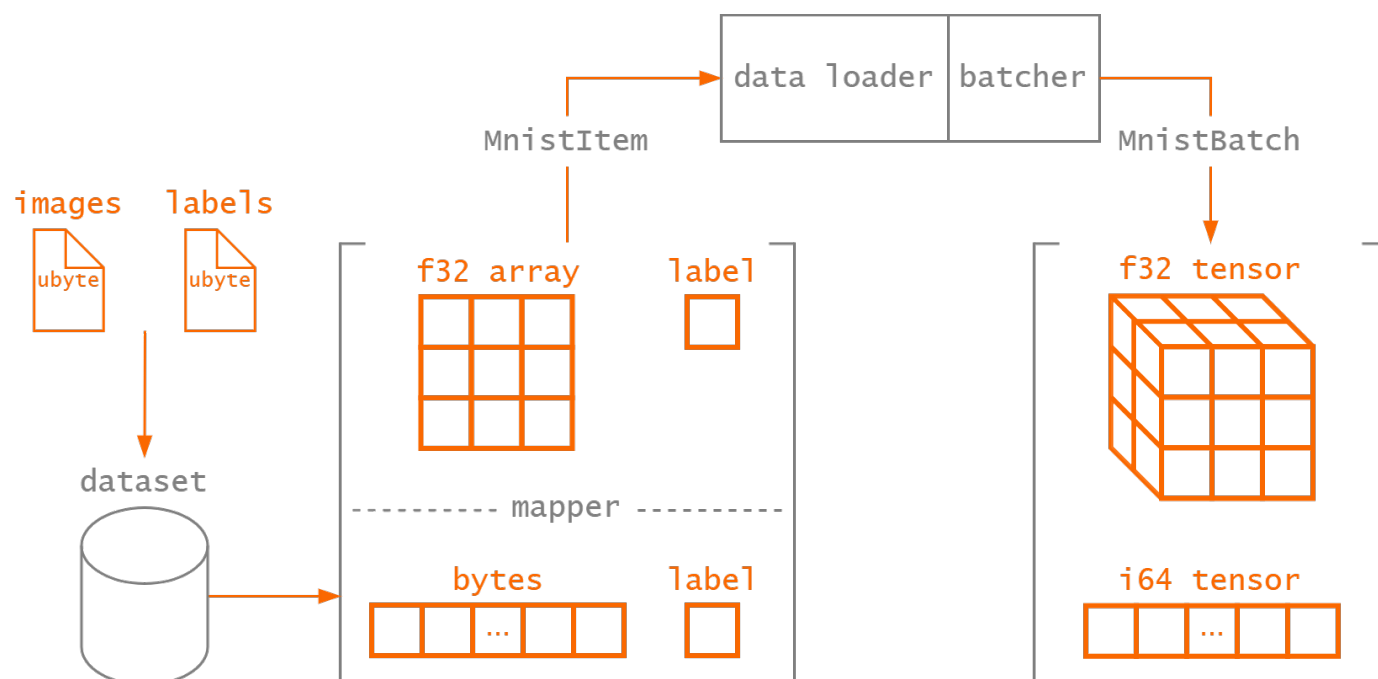
There is no streaming dataset API with Burn, and this is by design! The `learner struct` will iterate multiple times over the dataset and only checkpoint when done. You can consider the length of the dataset as the number of iterations before performing checkpointing and running the validation. There is nothing stopping you from returning different items even when called with the same `index` multiple times.

## How Is The Dataset Used?

During training, the dataset is used to access the data samples and, for most use cases in supervised learning, their corresponding ground-truth labels. Remember that the `Dataset` trait implementation is responsible to retrieve the data from its source, usually some sort of data storage. At this point, the dataset could be naively iterated over to provide the model a

single sample to process at a time, but this is not very efficient.

Instead, we collect multiple samples that the model can process as a *batch* to fully leverage modern hardware (e.g., GPUs - which have impressive parallel processing capabilities). Since each data sample in the dataset can be collected independently, the data loading is typically done in parallel to further speed things up. In this case, we parallelize the data loading using a multi-threaded `BatchDataLoader` to obtain a sequence of items from the `Dataset` implementation. Finally, the sequence of items is combined into a batched tensor that can be used as input to a model with the `Batcher` trait implementation. Other tensor operations can be performed during this step to prepare the batch data, as is done in the [basic workflow guide](#). The process is illustrated in the figure below for the MNIST dataset.



Although we have conveniently implemented the `MnistDataset` used in the guide, we'll go over its implementation to demonstrate how the `Dataset` and `Batcher` traits are used.

The `MNIST dataset` of handwritten digits has a training set of 60,000 examples and a test set of 10,000 examples. A single item in the dataset is represented by a  $28 \times 28$  pixels black-and-white image (stored as raw bytes) with its corresponding label (a digit between 0 and 9). This is defined by the `MnistItemRaw` struct.

```
struct MnistItemRaw {
  pub image_bytes: Vec<u8>,
  pub label: u8,
}
```

With single-channel images of such low resolution, the entire training and test sets can be loaded in memory at once. Therefore, we leverage the already existing `InMemDataset` to

retrieve the raw images and labels data. At this point, the image data is still just a bunch of bytes, but we want to retrieve the *structured* image data in its intended form. For that, we can define a `MapperDataset` that transforms the raw image bytes to a 2D array image (which we convert to float while we're at it).

```
const WIDTH: usize = 28;
const HEIGHT: usize = 28;

pub struct MnistItem {
    /// Image as a 2D array of floats.
    pub image: [[f32; WIDTH]; HEIGHT],

    /// Label of the image.
    pub label: u8,
}

struct BytesToImage;

impl Mapper<MnistItemRaw, MnistItem> for BytesToImage {
    /// Convert a raw MNIST item (image bytes) to a MNIST item (2D array image).
    fn map(&self, item: &MnistItemRaw) -> MnistItem {
        // Ensure the image dimensions are correct.
        debug_assert_eq!(item.image_bytes.len(), WIDTH * HEIGHT);

        // Convert the image to a 2D array of floats.
        let mut image_array = [[0f32; WIDTH]; HEIGHT];
        for (i, pixel) in item.image_bytes.iter().enumerate() {
            let x = i % WIDTH;
            let y = i / HEIGHT;
            image_array[y][x] = *pixel as f32;
        }

        MnistItem {
            image: image_array,
            label: item.label,
        }
    }
}

type MappedDataset = MapperDataset<InMemDataset<MnistItemRaw>, BytesToImage, MnistItemRaw>;

pub struct MnistDataset {
    dataset: MappedDataset,
}
```

To construct the `MnistDataset`, the data source must be parsed into the expected `MappedDataset` type. Since both the train and test sets use the same file format, we can separate the functionality to load the `train()` and `test()` dataset.

```

impl MnistDataset {
  /// Creates a new train dataset.
  pub fn train() -> Self {
    Self::new("train")
  }

  /// Creates a new test dataset.
  pub fn test() -> Self {
    Self::new("test")
  }

  fn new(split: &str) -> Self {
    // Download dataset
    let root = MnistDataset::download(split);

    // Parse data as vector of images bytes and vector of labels
    let images: Vec<Vec<u8>> = MnistDataset::read_images(&root, split);
    let labels: Vec<u8> = MnistDataset::read_labels(&root, split);

    // Collect as vector of MnistItemRaw
    let items: Vec<_> = images
      .into_iter()
      .zip(labels)
      .map(|(image_bytes, label)| MnistItemRaw { image_bytes, label })
      .collect();

    // Create the MapperDataset for InMemDataset<MnistItemRaw> to transform
    // items (MnistItemRaw -> MnistItem)
    let dataset = InMemDataset::new(items);
    let dataset = MapperDataset::new(dataset, BytesToImage);

    Self { dataset }
  }
}

```

Since the `MnistDataset` simply wraps a `MapperDataset` instance with `InMemDataset`, we can easily implement the `Dataset` trait.

```

impl Dataset<MnistItem> for MnistDataset {
  fn get(&self, index: usize) -> Option<MnistItem> {
    self.dataset.get(index)
  }

  fn len(&self) -> usize {
    self.dataset.len()
  }
}

```

The only thing missing now is the `Batcher`, which we already went over [in the basic workflow guide](#). The `Batcher` takes a list of `MnistItem` retrieved by the dataloader as input and returns a batch of images as a 3D tensor along with their targets.

# Custom Training Loops

Even though Burn comes with a project dedicated to simplifying training, it doesn't mean that you have to use it. Sometimes you may have special needs for your training, and it might be faster to just reimplement the training loop yourself. Also, you may just prefer implementing your own training loop instead of using a pre-built one in general.

Burn's got you covered!

We will start from the same example shown in the [basic workflow](#) section, but without using the `Learner` struct.

```

#[derive(Config)]
pub struct MnistTrainingConfig {
    #[config(default = 10)]
    pub num_epochs: usize,
    #[config(default = 64)]
    pub batch_size: usize,
    #[config(default = 4)]
    pub num_workers: usize,
    #[config(default = 42)]
    pub seed: u64,
    #[config(default = 1e-4)]
    pub lr: f64,
    pub model: ModelConfig,
    pub optimizer: AdamConfig,
}

pub fn run<B: AutodiffBackend>(device: &B::Device) {
    // Create the configuration.
    let config_model = ModelConfig::new(10, 1024);
    let config_optimizer = AdamConfig::new();
    let config = MnistTrainingConfig::new(config_model, config_optimizer);

    B::seed(config.seed);

    // Create the model and optimizer.
    let mut model = config.model.init(device);
    let mut optim = config.optimizer.init();

    // Create the batcher.
    let batcher_train = MnistBatcher::<B>::new(device.clone());
    let batcher_valid = MnistBatcher::<B::InnerBackend>::new(device.clone());

    // Create the dataloaders.
    let dataloader_train = DataLoaderBuilder::new(batcher_train)
        .batch_size(config.batch_size)
        .shuffle(config.seed)
        .num_workers(config.num_workers)
        .build(MnistDataset::train());

    let dataloader_test = DataLoaderBuilder::new(batcher_valid)
        .batch_size(config.batch_size)
        .shuffle(config.seed)
        .num_workers(config.num_workers)
        .build(MnistDataset::test());

    ...
}

```

As seen with the previous example, setting up the configurations and the dataloader hasn't changed. Now, let's move forward and write our own training loop:



```

pub fn run<B: AutodiffBackend>(device: B::Device) {
    ...

    // Iterate over our training and validation loop for X epochs.
    for epoch in 1..config.num_epochs + 1 {
        // Implement our training loop.
        for (iteration, batch) in dataloader_train.iter().enumerate() {
            let output = model.forward(batch.images);
            let loss = CrossEntropyLoss::new(None, &output.device())
                .forward(output.clone(), batch.targets.clone());
            let accuracy = accuracy(output, batch.targets);

            println!(
                "[Train - Epoch {} - Iteration {}] Loss {:.3} | Accuracy {:.3}
%",
                epoch,
                iteration,
                loss.clone().into_scalar(),
                accuracy,
            );

            // Gradients for the current backward pass
            let grads = loss.backward();
            // Gradients linked to each parameter of the model.
            let grads = GradientsParams::from_grads(grads, &model);
            // Update the model using the optimizer.
            model = optim.step(config.lr, model, grads);
        }

        // Get the model without autodiff.
        let model_valid = model.valid();

        // Implement our validation loop.
        for (iteration, batch) in dataloader_test.iter().enumerate() {
            let output = model_valid.forward(batch.images);
            let loss = CrossEntropyLoss::new(None, &output.device())
                .forward(output.clone(), batch.targets.clone());
            let accuracy = accuracy(output, batch.targets);

            println!(
                "[Valid - Epoch {} - Iteration {}] Loss {} | Accuracy {}",
                epoch,
                iteration,
                loss.clone().into_scalar(),
                accuracy,
            );
        }
    }
}

```

In the previous code snippet, we can observe that the loop starts from epoch 1 and goes up

to `num_epochs` . Within each epoch, we iterate over the training dataloader. During this process, we execute the forward pass, which is necessary for computing both the loss and accuracy. To maintain simplicity, we print the results to `stdout`.

Upon obtaining the loss, we can invoke the `backward()` function, which returns the gradients specific to each variable. It's important to note that we need to map these gradients to their corresponding parameters using the `GradientsParams` type. This step is essential because you might run multiple different autodiff graphs and accumulate gradients for each parameter id.

Finally, we can perform the optimization step using the learning rate, the model, and the computed gradients. It's worth mentioning that, unlike PyTorch, there's no need to register the gradients with the optimizer, nor do you have to call `zero_grad` . The gradients are automatically consumed during the optimization step. If you're interested in gradient accumulation, you can easily achieve this by using the `GradientsAccumulator` .

```
let mut accumulator = GradientsAccumulator::new();
let grads = model.backward();
let grads = GradientsParams::from_grads(grads, &model);
accumulator.accumulate(&model, grads); ...
let grads = accumulator.grads(); // Pop the accumulated gradients.
```

Note that after each epoch, we include a validation loop to assess our model's performance on previously unseen data. To disable gradient tracking during this validation step, we can invoke `model.valid()` , which provides a model on the inner backend without autodiff capabilities. It's important to emphasize that we've declared our validation batcher to be on the inner backend, specifically `MnistBatcher<B::InnerBackend>` ; not using `model.valid()` will result in a compilation error.

You can find the code above available as an [example](#) for you to test.

## Multiple optimizers

It's common practice to set different learning rates, optimizer parameters, or use different optimizers entirely, for different parts of a model. In Burn, each `GradientParams` can contain only a subset of gradients to actually apply with an optimizer. This allows you to flexibly mix and match optimizers!

```
// Start with calculating all gradients
let grads = loss.backward();

// Now split the gradients into various parts.
let grads_conv1 = GradientParams::from_module(&mut grads, &model.conv1);
let grads_conv2 = GradientParams::from_module(&mut grads, &model.conv2);

// You can step the model with these gradients, using different learning
// rates for each param. You could also use an entirely different optimizer
// here!
model = optim.step(config.lr * 2.0, model, grads_conv1);
model = optim.step(config.lr * 4.0, model, grads_conv2);

// For even more granular control you can split off individual parameter
// eg. a linear bias usually needs a smaller learning rate.
if let Some(bias) == model.linear1.bias {
    let grads_bias = GradientParams::from_params(&mut grads, &model.linear1,
    &[bias.id]);
    model = optim.step(config.lr * 0.1, model, grads_bias);
}

// Note that above calls remove gradients, so we can just get all "remaining"
// gradients.
let grads = GradientsParams::from_grads(grads, &model);
model = optim.step(config.lr, model, grads);
```

## Custom Type

The explanations above demonstrate how to create a basic training loop. However, you may find it beneficial to organize your program using intermediary types. There are various ways to do this, but it requires getting comfortable with generics.

If you wish to group the optimizer and the model into the same structure, you have several options. It's important to note that the optimizer trait depends on both the `AutodiffModule` trait and the `AutodiffBackend` trait, while the module only depends on the `AutodiffBackend` trait.

Here's a closer look at how you can create your types:

**Create a struct that is generic over the backend and the optimizer, with a predefined model.**

```
struct Learner<B, O>
where
  B: AutodiffBackend,
{
  model: Model<B>,
  optim: O,
}
```

This is quite straightforward. You can be generic over the backend since it's used with the concrete type `Model` in this case.

### Create a struct that is generic over the model and the optimizer.

```
struct Learner<M, O> {
  model: M,
  optim: O,
}
```

This option is a quite intuitive way to declare the struct. You don't need to write type constraints with a `where` statement when defining a struct; you can wait until you implement the actual function. However, with this struct, you may encounter some issues when trying to implement code blocks to your struct.

```
impl<B, M, O> Learner<M, O>
where
  B: AutodiffBackend,
  M: AutodiffModule<B>,
  O: Optimizer<M, B>,
{
  pub fn step(&mut self, _batch: MnistBatch<B>) {
    //
  }
}
```

This will result in the following compilation error:

1. the type parameter ``B`` is not constrained by the `impl` trait, `self` type, or predicates  
unconstrained type parameter [E0207]

To resolve this issue, you have two options. The first one is to make your function generic over the backend and add your trait constraint within its definition:

```
#[allow(dead_code)]
impl<M, O> Learner2<M, O> {
    pub fn step<B: AutodiffBackend>(&mut self, _batch: MnistBatch<B>)
    where
        B: AutodiffBackend,
        M: AutodiffModule<B>,
        O: Optimizer<M, B>,
    {
        //
    }
}
```

However, some people may prefer to have the constraints on the implementation block itself. In that case, you can make your struct generic over the backend using `PhantomData<B>`.

**Create a struct that is generic over the backend, the model, and the optimizer.**

```
struct Learner3<B, M, O> {
    model: M,
    optim: O,
    _b: PhantomData<B>,
}
```

You might wonder why `PhantomData` is required. Each generic argument must be used as a field when declaring a struct. When you don't need the generic argument, you can use `PhantomData` to mark it as a zero sized type.

These are just some suggestions on how to define your own types, but you are free to use any pattern that you prefer.

# Saving and Loading Models

Saving your trained machine learning model is quite easy, no matter the output format you choose. As mentioned in the [Record](#) section, different formats are supported to serialize/deserialize models. By default, we use the `NamedMpkFileRecorder` which uses the [MessagePack](#) binary serialization format with the help of [smp\\_serde](#).

```
// Save model in MessagePack format with full precision
let recorder = NamedMpkFileRecorder::<FullPrecisionSettings>::new();
model
    .save_file(model_path, &recorder)
    .expect("Should be able to save the model");
```

Note that the file extension is automatically handled by the recorder depending on the one you choose. Therefore, only the file path and base name should be provided.

Now that you have a trained model saved to your disk, you can easily load it in a similar fashion.

```
// Load model in full precision from MessagePack file
let recorder = NamedMpkFileRecorder::<FullPrecisionSettings>::new();
model
    .load_file(model_path, &recorder, device)
    .expect("Should be able to load the model weights from the provided file");
```

**Note:** models can be saved in different output formats, just make sure you are using the correct recorder type when loading the saved model. Type conversion between different precision settings is automatically handled, but formats are not interchangeable. A model can be loaded from one format and saved to another format, just as long as you load it back with the new recorder type afterwards.

## Initialization from Recorded Weights

The most straightforward way to load weights for a module is simply by using the generated method [load\\_record](#). Note that parameter initialization is lazy, therefore no actual tensor allocation and GPU/CPU kernels are executed before the module is used. This means that you can use `init(device)` followed by `load_record(record)` without any meaningful performance cost.

```
// Create a dummy initialized model to save
let device = Default::default();
let model = Model::<MyBackend>::init(&device);

// Save model in MessagePack format with full precision
let recorder = NamedMpkFileRecorder::<FullPrecisionSettings>::new();
model
    .save_file(model_path, &recorder)
    .expect("Should be able to save the model");
```

Afterwards, the model can just as easily be loaded from the record saved on disk.

```
// Load model record on the backend's default device
let record: ModelRecord<MyBackend> =
    NamedMpkFileRecorder::<FullPrecisionSettings>::new()
        .load(model_path.into(), &device)
        .expect("Should be able to load the model weights from the provided file");

// Initialize a new model with the loaded record/weights
let model = Model::init(&device).load_record(record);
```

## No Storage, No Problem!

For applications where file storage may not be available (or desired) at runtime, you can use the `BinBytesRecorder`.

In the previous examples we used a `FileRecorder` based on the MessagePack format, which could be replaced with [another file recorder](#) of your choice. To embed a model as part of your runtime application, first save the model to a binary file with `BinFileRecorder`.

```
// Save model in binary format with full precision
let recorder = BinFileRecorder::<FullPrecisionSettings>::new();
model
    .save_file(model_path, &recorder)
    .expect("Should be able to save the model");
```

Then, in your final application, include the model and use the `BinBytesRecorder` to load it.

Embedding the model as part of your application is especially useful for smaller models but not recommended for very large models as it would significantly increase the binary size as well as consume a lot more memory at runtime.

```
// Include the model file as a reference to a byte array
static MODEL_BYTES: &[u8] = include_bytes!("path/to/model.bin");

// Load model binary record in full precision
let record = BinBytesRecorder::<FullPrecisionSettings>::default()
    .load(MODEL_BYTES.to_vec(), device)
    .expect("Should be able to load model the model weights from bytes");

// Load that record with the model
model.load_record(record);
```

This example assumes that the model was already created before loading the model record. If instead you want to skip the random initialization and directly initialize the weights with the provided record, you could adapt this like the [previous example](#).



# Importing Models

The Burn project supports the import of models from various frameworks, emphasizing efficiency and compatibility. Currently, it handles two primary model formats:

1. [ONNX](#): Facilitates direct import, ensuring the model's performance and structure are maintained.
2. [PyTorch](#): Enables the loading of PyTorch model weights into Burn's native model architecture, ensuring seamless integration.

# Importing ONNX Models in Burn

## Table of Contents

1. [Introduction](#)
2. [Why Import Models?](#)
3. [Understanding ONNX](#)
4. [Burn's ONNX Support](#)
5. [Step-by-Step Guide](#)
6. [Advanced Configuration](#)
7. [Loading and Using Models](#)
8. [Troubleshooting](#)
9. [Examples and Resources](#)
10. [Conclusion](#)

## Introduction

As the field of deep learning continues to evolve, the need for interoperability between different frameworks becomes increasingly important. Burn, a modern deep learning framework in Rust, recognizes this need and provides robust support for importing models from other popular frameworks. This section focuses on importing [ONNX \(Open Neural Network Exchange\)](#) models into Burn, enabling you to leverage pre-trained models and seamlessly integrate them into your Rust-based deep learning projects.

## Why Import Models?

Importing pre-trained models offers several advantages:

1. **Time-saving:** Avoid the need to train models from scratch, which can be time-consuming and resource-intensive.
2. **Access to state-of-the-art architectures:** Utilize cutting-edge models developed by researchers and industry leaders.
3. **Transfer learning:** Fine-tune imported models for your specific tasks, benefiting from knowledge transfer.
4. **Consistency across frameworks:** Ensure consistent performance when moving from

one framework to another.

## Understanding ONNX

ONNX (Open Neural Network Exchange) is an open format designed to represent machine learning models. Key features include:

- **Framework agnostic:** ONNX provides a common format that works across various deep learning frameworks.
- **Comprehensive representation:** It captures both the model architecture and trained weights.
- **Wide support:** Many popular frameworks like PyTorch, TensorFlow, and scikit-learn support ONNX export.

By using ONNX, you can easily move models between different frameworks and deployment environments.

## Burn's ONNX Support

Burn takes a unique approach to ONNX import, offering several advantages:

1. **Native Rust code generation:** ONNX models are translated into Rust source code, allowing for deep integration with Burn's ecosystem.
2. **Compile-time optimization:** The generated Rust code can be optimized by the Rust compiler, potentially improving performance.
3. **No runtime dependency:** Unlike some solutions that require an ONNX runtime, Burn's approach eliminates this dependency.
4. **Trainability:** Imported models can be further trained or fine-tuned using Burn.
5. **Portability:** The generated Rust code can be compiled for various targets, including WebAssembly and embedded devices.
6. **Any Burn Backend:** The imported models can be used with any of Burn's backends.

## Step-by-Step Guide

Let's walk through the process of importing an ONNX model into a Burn project:

## Step 1: Update `build.rs`

First, add the `burn-import` crate to your `Cargo.toml`:

```
[build-dependencies]
burn-import = "~0.15"
```

Then, in your `build.rs` file:

```
use burn_import::onnx::ModelGen;

fn main() {
    ModelGen::new()
        .input("src/model/my_model.onnx")
        .out_dir("model/")
        .run_from_script();
}
```

This script uses `ModelGen` to generate Rust code from your ONNX model during the build process.

## Step 2: Modify `mod.rs`

In your `src/model/mod.rs` file, include the generated code:

```
pub mod my_model {
    include!(concat!(env!("OUT_DIR"), "/model/my_model.rs"));
}
```

This makes the generated model code available in your project.

## Step 3: Use the Imported Model

Now you can use the imported model in your Rust code:

```

use burn::tensor;
use burn_ndarray::{NdArray, NdArrayDevice};
use model::my_model::Model;

fn main() {
    let device = NdArrayDevice::default();

    // Create model instance and load weights from target dir default device.
    // (see more load options below in "Loading and Using Models" section)
    let model: Model<NdArray<f32>> = Model::default();

    // Create input tensor (replace with your actual input)
    let input = tensor::Tensor::<NdArray<f32>, 4>::zeros([1, 3, 224, 224],
&device);

    // Perform inference
    let output = model.forward(input);

    println!("Model output: {:?}", output);
}

```

## Advanced Configuration

The `ModelGen` struct offers several configuration options:

```

ModelGen::new()
    .input("path/to/model.onnx")
    .out_dir("model/")
    .record_type(RecordType::NamedMpk)
    .half_precision(false)
    .embed_states(false)
    .run_from_script();

```

- `record_type` : Specifies the format for storing weights (Bincode, NamedMpk, NamedMpkGz, or PrettyJson).
- `half_precision` : Use half-precision (f16) for weights to reduce model size.
- `embed_states` : Embed model weights directly in the generated Rust code. Note: This requires record type `Bincode`.

## Loading and Using Models

Depending on your configuration, you can load models in different ways:

```
// Create a new model instance with device. Initializes weights randomly and lazily.
// You can load weights via `load_record` afterwards.
let model = Model::::new(&device);

// Load from a file (must specify weights file in the target output directory or copy it from there).
// File type should match the record type specified in `ModelGen`.
let model = Model::::from_file("path/to/weights", &device);

// Load from embedded weights (if embed_states was true)
let model = Model::::from_embedded();

// Load from the out director location and load to default device (useful for testing)
let model = Model::::default();
```

## Troubleshooting

Here are some common issues and their solutions:

1. **Unsupported ONNX operator:** If you encounter an error about an unsupported operator, check the [list of supported ONNX operators](#). You may need to simplify your model or wait for support to be added.
2. **Build errors:** Ensure that your `burn-import` version matches your Burn version. Also, check that the ONNX file path in `build.rs` is correct.
3. **Runtime errors:** If you get errors when running your model, double-check that your input tensors match the expected shape and data type of your model.
4. **Performance issues:** If your imported model is slower than expected, try using the `half_precision` option to reduce memory usage, or experiment with different `record_type` options.
5. **Artifact Files:** You can view the generated Rust code and weights files in the `OUT_DIR` directory specified in `build.rs` (usually `target/debug/build/<project>/out`).

## Examples and Resources

For more detailed examples, check out:

1. [MNIST Inference Example](#)
2. [SqueezeNet Image Classification](#)

These examples demonstrate real-world usage of ONNX import in Burn projects.

## Conclusion

Importing ONNX models into Burn opens up a world of possibilities, allowing you to leverage pre-trained models from other frameworks while taking advantage of Burn's performance and Rust's safety features. By following this guide, you should be able to seamlessly integrate ONNX models into your Burn projects, whether for inference, fine-tuning, or as a starting point for further development.

Remember that the `burn-import` crate is actively developed, with ongoing work to support more ONNX operators and improve performance. Stay tuned to the Burn repository for updates and new features!



**Note:** The `burn-import` crate is in active development. For the most up-to-date information on supported ONNX operators, please refer to the [official documentation](#).

---

# PyTorch Model

## Introduction

Whether you've trained your model in PyTorch or you want to use a pre-trained model from PyTorch, you can import them into Burn. Burn supports importing PyTorch model weights with `.pt` file extension. Compared to ONNX models, `.pt` files only contain the weights of the model, so you will need to reconstruct the model architecture in Burn.

Here in this document we will show the full workflow of exporting a PyTorch model and importing it. Also you can refer to this [Transitioning From PyTorch to Burn](#) tutorial on importing a more complex model.

## How to export a PyTorch model

If you have a PyTorch model that you want to import into Burn, you will need to export it first, unless you are using a pre-trained published model. To export a PyTorch model, you can use the `torch.save` function.

The following is an example of how to export a PyTorch model:

```
import torch
import torch.nn as nn

class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.conv1 = nn.Conv2d(2, 2, (2,2))
        self.conv2 = nn.Conv2d(2, 2, (2,2), bias=False)

    def forward(self, x):
        x = self.conv1(x)
        x = self.conv2(x)
        return x

if __name__ == "__main__":
    torch.manual_seed(42) # To make it reproducible
    model = Net().to(torch.device("cpu"))
    model_weights = model.state_dict()
    torch.save(model_weights, "conv2d.pt")
```



Use [Netron](#) to view the exported model. You should see something like this:

conv2	
weight	$\langle 2 \times 2 \times 2 \times 2 \rangle$

conv1	
weight	$\langle 2 \times 2 \times 2 \times 2 \rangle$
bias	$\langle 2 \rangle$

## How to import a PyTorch model

1. Define the model in Burn:

```

use burn::{
    nn::conv::{Conv2d, Conv2dConfig},
    prelude::*,
};

#[derive(Module, Debug)]
pub struct Net<B: Backend> {
    conv1: Conv2d<B>,
    conv2: Conv2d<B>,
}

impl<B: Backend> Net<B> {
    /// Create a new model.
    pub fn init(device: &B::Device) -> Self {
        let conv1 = Conv2dConfig::new([2, 2], [2, 2])
            .init(device);
        let conv2 = Conv2dConfig::new([2, 2], [2, 2])
            .with_bias(false)
            .init(device);
        Self { conv1, conv2 }
    }

    /// Forward pass of the model.
    pub fn forward(&self, x: Tensor<B, 4>) -> Tensor<B, 4> {
        let x = self.conv1.forward(x);
        self.conv2.forward(x)
    }
}

```

2. Load the model weights from the exported PyTorch model (2 options):

a) *Dynamically*, but this requires burn-import runtime dependency:

```
use crate::model;

use burn::record::{FullPrecisionSettings, Recorder};
use burn_import::pytorch::PyTorchFileRecorder;

type Backend = burn_ndarray::NdArray<f32>;

fn main() {
    let device = Default::default();
    let record = PyTorchFileRecorder::<FullPrecisionSettings>::default()
        .load("./conv2d.pt".into(), &device)
        .expect("Should decode state successfully");

    let model = model::Net::<Backend>::init(&device).load_record(record);
}
```

b) *Pre-converted* to Burn's binary format:

```

// Convert the PyTorch model to Burn's binary format in
// build.rs or in a separate executable. Then, include the generated file
// in your project. See `examples/pytorch-import` for an example.

use crate::model;

use burn::record::{FullPrecisionSettings, NamedMpkFileRecorder, Recorder};
use burn_import::pytorch::PyTorchFileRecorder;

type Backend = burn_ndarray::NdArray<f32>;

fn main() {
    let device = Default::default();
    let recorder = PyTorchFileRecorder::<FullPrecisionSettings>::default()
    let record: model::NetRecord<B> = recorder
        .load("./conv2d.pt".into(), &device)
        .expect("Should decode state successfully");

    // Save the model record to a file.
    let recorder =
NamedMpkFileRecorder::<FullPrecisionSettings>::default();
    recorder
        .record(record, "MY_FILE_OUTPUT_PATH".into())
        .expect("Failed to save model record");
}

/// Load the model from the file in your source code (not in build.rs or
script).
fn load_model() -> Net::<Backend> {
    let device = Default::default();
    let record = NamedMpkFileRecorder::<FullPrecisionSettings>::default()
        .load("./MY_FILE_OUTPUT_PATH".into(), &device)
        .expect("Should decode state successfully");

    Net::<Backend>::init(&device).load_record(record)
}

```

## Extract Configuration

In some cases, models may require additional configuration settings, which are often included in a `.pt` file during export. The `config_from_file` function from the `burn-import` cargo package allows for the extraction of these configurations directly from the `.pt` file. The extracted configuration can then be used to initialize the model in Burn. Here is an example of how to extract the configuration from a `.pt` file:

```
use std::collections::HashMap;

use burn::config::Config;
use burn_import::pytorch::config_from_file;

#[derive(Debug, Config)]
struct NetConfig {
    n_head: usize,
    n_layer: usize,
    d_model: usize,
    // Candle's pickle has a bug with float serialization
    // https://github.com/huggingface/candle/issues/1729
    // some_float: f64,
    some_int: i32,
    some_bool: bool,
    some_str: String,
    some_list_int: Vec<i32>,
    some_list_str: Vec<String>,
    // Candle's pickle has a bug with float serialization
    // https://github.com/huggingface/candle/issues/1729
    // some_list_float: Vec<f64>,
    some_dict: HashMap<String, String>,
}

fn main() {
    let path = "weights_with_config.pt";
    let top_level_key = Some("my_config");
    let config: NetConfig = config_from_file(path, top_level_key).unwrap();
    println!("{:?}", config);

    // After extracting, it's recommended you save it as a json file.
    config.save("my_config.json").unwrap();
}
```

## Troubleshooting

### Adjusting the source model architecture

If your target model differs structurally from the model you exported, `PyTorchFileRecorder`

allows you to change the attribute names and the order of the attributes. For example, if you exported a model with the following structure:

```
class ConvModule(nn.Module):
    def __init__(self):
        super(ConvModule, self).__init__()
        self.conv1 = nn.Conv2d(2, 2, (2,2))
        self.conv2 = nn.Conv2d(2, 2, (2,2), bias=False)

    def forward(self, x):
        x = self.conv1(x)
        x = self.conv2(x)
        return x

class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.conv = ConvModule()

    def forward(self, x):
        x = self.conv(x)
        return x
```

But you need to import it into a model with the following structure:

```
#[derive(Module, Debug)]
pub struct Net<B: Backend> {
    conv1: Conv2d<B>,
    conv2: Conv2d<B>,
}
```

Which produces the following weights structure (viewed in [Netron](#)):

conv.conv2

weight <2×2×2×2>

conv.conv1

weight <2×2×2×2>

bias <2>

You can use the `PyTorchFileRecorder` to change the attribute names and the order of the attributes by specifying a regular expression (See [regex::Regex::replace](#) and [try it online](#)) to match the attribute name and a replacement string in `LoadArgs` :

```

let device = Default::default();
let load_args = LoadArgs::new("tests/key_remap/key_remap.pt".into())
    // Remove "conv" prefix, e.g. "conv.conv1" -> "conv1"
    .with_key_remap("conv\\.(.*)", "$1");

let record = PyTorchFileRecorder::<FullPrecisionSettings>::default()
    .load(load_args, &device)
    .expect("Should decode state successfully");

let model = Net::<Backend>::init(&device).load_record(record);

```

## Printing the source model keys and tensor information

If you are unsure about the keys in the source model, you can print them using the following code:

```

let device = Default::default();
let load_args = LoadArgs::new("tests/key_remap/key_remap.pt".into())
    // Remove "conv" prefix, e.g. "conv.conv1" -> "conv1"
    .with_key_remap("conv\\.(.*)", "$1")
    .with_debug_print(); // Print the keys and remapped keys

let record = PyTorchFileRecorder::<FullPrecisionSettings>::default()
    .load(load_args, &device)
    .expect("Should decode state successfully");

let model = Net::<Backend>::init(&device).load_record(record);

```

Here is an example of the output:

Debug information of keys and tensor shapes:

---

Original Key: conv.conv1.bias

Remapped Key: conv1.bias

Shape: [2]

Dtype: F32

---

Original Key: conv.conv1.weight

Remapped Key: conv1.weight

Shape: [2, 2, 2, 2]

Dtype: F32

---

Original Key: conv.conv2.weight

Remapped Key: conv2.weight

Shape: [2, 2, 2, 2]

Dtype: F32

---

## Non-contiguous indices in the source model

Sometimes the indices of the source model are non-contiguous. For example, the source model has:

```
"model.ax.0.weight",
"model.ax.0.bias",
"model.ax.2.weight",
"model.ax.2.bias",
"model.ax.4.weight",
"model.ax.4.bias",
"model.ax.6.weight",
"model.ax.6.bias",
"model.ax.8.weight",
"model.ax.8.bias",
```

This may occur when `model.ax` attribute (in the above example) uses `Sequential` to define the layers and the skipped items do not have weight tensors, such as a `ReLU` layer. PyTorch simply skips the layers without weight tensors, resulting in non-contiguous indices. In this case, `PyTorchFileRecorder` automatically corrects the indices to be contiguous preserving the order of the weights resulting in:

```
"model.ax.0.weight",
"model.ax.0.bias",
"model.ax.1.weight",
"model.ax.1.bias",
"model.ax.2.weight",
"model.ax.2.bias",
"model.ax.3.weight",
"model.ax.3.bias",
"model.ax.4.weight",
"model.ax.4.bias",
```

## Loading the model weights to a partial model

`PyTorchFileRecorder` enables selective weight loading into partial models. For instance, in a model with both an encoder and a decoder, it's possible to load only the encoder weights. This is done by defining the encoder in Burn, allowing the loading of its weights while excluding the decoder's.

## Specifying the top-level key for `state_dict`



Sometimes the `state_dict` is nested under a top-level key along with other metadata as in a [general checkpoint](#). For example, the `state_dict` of the whisper model is nested under `model_state_dict` key. In this case, you can specify the top-level key in `LoadArgs` :

```
let device = Default::default();
let load_args = LoadArgs::new("tiny.en.pt".into())
    .with_top_level_key("my_state_dict");

let record = PyTorchFileRecorder::<FullPrecisionSettings>::default()
    .load(load_args, &device)
    .expect("Should decode state successfully")
```

## Models containing enum modules

Burn supports models containing enum modules with new-type variants (tuple with one item). Importing weights for such models is automatically supported by the `PyTorchFileRecorder`. However, it should be noted that since the source weights file does not contain the enum variant information, the enum variant is picked based on the enum variant type. Let's consider the following example:

```
#[derive(Module, Debug)]
pub enum Conv<B: Backend> {
    DwsConv(DwsConv<B>),
    Conv(Conv2d<B>),
}

#[derive(Module, Debug)]
pub struct DwsConv<B: Backend> {
    dconv: Conv2d<B>,
    pconv: Conv2d<B>,
}

#[derive(Module, Debug)]
pub struct Net<B: Backend> {
    conv: Conv<B>,
}
```

If the source weights file contains weights for `DwsConv` , such as the following keys:

```
---
Key: conv.dconv.bias
Shape: [2]
Dtype: F32
---
Key: conv.dconv.weight
Shape: [2, 1, 3, 3]
Dtype: F32
---
Key: conv.pconv.bias
Shape: [2]
Dtype: F32
---
Key: conv.pconv.weight
Shape: [2, 2, 1, 1]
Dtype: F32
```

The weights will be imported into the `DwsConv` variant of the `Conv` enum module.

If the variant types are identical, then the first variant is picked. Generally, it won't be a problem since the variant types are usually different.

## Current known issues

1. [Candle's pickle does not currently unpack boolean tensors.](#)

# Models and Pre-Trained Weights

The [models](#) repository contains definitions of different deep learning models with examples for different domains like computer vision and natural language processing.

This includes image classification models such as [MobileNetV2](#) , [SqueezeNet](#) and [ResNet](#) , object detection models such as [YOLOX](#) and language models like [BERT](#) and [RoBERTa](#) .

Be sure to check out the up-to-date [collection of models](#) to get you started. Pre-trained weights are available for every supported architecture in this collection. You will also find a spotlight of [community contributed models](#).

# Quantization (Beta)

Quantization techniques perform computations and store tensors in lower precision data types like 8-bit integer instead of floating point precision. There are multiple approaches to quantize a deep learning model categorized as:

- Post-training quantization (PTQ)
- Quantization aware training (QAT)

In post-training quantization, the model is trained in floating point precision and later converted to the lower precision data type.

There are two types of post-training quantization:

1. Static quantization: quantizes the weights and activations of the model. Quantizing the activations statically requires data to be calibrated (i.e., recording the activation values to compute the optimal quantization parameters with representative data).
2. Dynamic quantization: quantized the weights ahead of time (like static quantization) but the activations are dynamically at runtime.

Sometimes post-training quantization is not able to achieve acceptable task accuracy. This is where quantization aware training comes into play, as it models the effects of quantization during training. Quantization errors are thus modeled in the forward and backward passes using fake quantization modules, which helps the model learn representations that are more robust to the reduction in precision.

Quantization support in Burn is currently in active development.

It supports the following modes on some backends:

- Static per-tensor quantization to signed 8-bit integer ( `int8` )

No integer operations are currently supported, which means tensors are dequantized to perform the operations in floating point precision.

## Module Quantization

Quantizing the weights of your model after training is quite simple. We have access to the weight tensors and can collect their statistics, such as the min and max value when using `MinMaxCalibration`, to compute the quantization parameters.

```
// Quantization config
let mut quantizer = Quantizer {
    calibration: MinMaxCalibration {},
    scheme: QuantizationScheme::PerTensorSymmetric(QuantizationType::QInt8),
};

// Quantize the weights
let model = model.quantize_weights(&mut quantizer);
```

---

Given that all operations are currently performed in floating point precision, it might be wise to dequantize the module parameters before inference. This allows us to save disk space by storing the model in reduced precision while preserving the inference speed.

This can easily be implemented with a `ModuleMapper`.

```
/// Module mapper used to dequantize the model params being loaded.
pub struct Dequantize {}

impl<B: Backend> ModuleMapper<B> for Dequantize {
    fn map_float<const D: usize>(
        &mut self,
        _id: ParamId,
        tensor: Tensor<B, D>,
    ) -> Tensor<B, D> {
        tensor.dequantize()
    }
}

// Load saved quantized model in floating point precision
model = model
    .load_file(file_path, recorder, &device)
    .expect("Should be able to load the quantized model weights")
    .map(&mut Dequantize {});
```

---

## Calibration

Calibration is the step during quantization where the range of all floating-point tensors is computed. This is pretty straightforward for weights since the actual range is known at *quantization-time* (weights are static), but activations require more attention.

To compute the quantization parameters, Burn supports the following `Calibration` methods.

Method	Description
MinMaxCalibration	Computes the quantization range mapping based on the running min and max values.

## Quantization Scheme

A quantization scheme defines the quantized type, quantization granularity and range mapping technique.

Burn currently supports the following `QuantizationType` variants.

Type	Description
QInt8	8-bit signed integer quantization.

Quantization parameters are defined based on the range of values to represent and can typically be calculated for the layer's entire weight tensor with per-tensor quantization or separately for each channel with per-channel quantization (commonly used with CNNs).

Burn currently supports the following `QuantizationScheme` variants.

Variant	Description
PerTensorAffine	Computes the quantization parameters for the whole tensor and applies an affine range mapping with zero point.
PerTensorSymmetric	Computes the quantization parameters for the whole tensor and applies a scale range mapping centered around 0.

# Advanced

In this section, we will go into advanced topics that extend beyond basic usage. Given Burn's exceptional flexibility, a lot of advanced use cases become possible.

Before going through this section, we strongly recommend exploring the [basic workflow](#) section and the [building blocks](#) section. Establishing a solid understanding of how the framework operates is crucial to comprehending the advanced concepts presented here. While you have the freedom to explore the advanced sections in any order you prefer, it's important to note that this section is not intended to be linear, contrary to preceding sections. Instead, it serves as a repository of use cases that you can refer to for guidance as needed.

# Backend Extension

Burn aims to be the most flexible deep learning framework. While it's crucial to maintain compatibility with a wide variety of backends, Burn provides the ability to extend the functionality of a backend implementation to suit your modeling requirements. This versatility is advantageous in numerous ways, such as supporting custom operations like flash attention or manually fusing operations for enhanced performance.

In this section, we will go into the process of extending a backend, providing multiple examples. But before we proceed, let's establish the fundamental principles that will empower you to craft your own backend extensions.

As you can observe, most types in Burn are generic over the Backend trait. This might give the impression that Burn operates at a high level over the backend layer. However, making the trait explicit instead of being chosen via a compilation flag was a thoughtful design decision. This explicitness does not imply that all backends must be identical; rather, it offers a great deal of flexibility when composing backends. The autodifferentiation backend trait (see [autodiff section](#)) is an example of how the backend trait has been extended to enable gradient computation with backpropagation. Furthermore, this design allows you to create your own backend extension. To achieve this, you need to design your own backend trait specifying which functions should be supported.

```
pub trait Backend: burn::tensor::backend::Backend {
    fn my_new_function(tensor: B::TensorPrimitive<2>) -> B::TensorPrimitive<2>
{
    // You can define a basic implementation reusing the Burn Backend API.
    // This can be useful since all backends will now automatically support
    // your model. But performance can be improved for this new
    // operation by implementing this block in specific backends.
}
}
```

You can then implement your new custom backend trait for any backend that you want to support:

```
impl<E: TchElement> Backend for burn_tch::LibTorch<E> {
    fn my_new_function(tensor: TchTensor<E, 2>) -> TchTensor<E, 2> {
        // My Tch implementation
    }
}

impl<E: NdArrayElement> Backend for burn_ndarray::NdArray<E> {
    // No specific implementation, but the backend can still be used.
}
```



You can support the backward pass using the same pattern.

```
impl<B: Backend> Backend for burn_autodiff::Autodiff<B> {
    // No specific implementation; autodiff will work with the default
    // implementation. Useful if you still want to train your model, but
    // observe performance gains mostly during inference.
}

impl<B: Backend> Backend for burn_autodiff::Autodiff<B> {
    fn my_new_function(tensor: AutodiffTensor<E, 2>) -> AutodiffTensor<E, 2> {
        // My own backward implementation, generic over my custom Backend trait.
        //
        // You can add a new method `my_new_function_backward` to your custom
        backend
        // trait if you want to invoke a custom kernel during the backward pass.
    }
}

impl<E: TchElement> Backend for burn_autodiff::Autodiff<burn_tch::LibTorch<E>>
{
    fn my_new_function(tensor: AutodiffTensor<E, 2>) -> AutodiffTensor<E, 2> {
        // My own backward implementation, generic over a backend implementation.
        //
        // This is another way to call a custom kernel for the backward pass that
        // doesn't require the addition of a new `backward` function in the
        custom backend.
        // This is useful if you don't want all backends to support training,
        reducing
        // the need for extra code when you know your model will only be trained
        on one
        // specific backend.
    }
}
```

The specifics of each implementation will be covered by the examples provided in this section. The `cubec1` compiler frontend is the recommended method of implementing custom kernels, since it supports multiple backends, including `wgpu` and `CUDA`, and is the way first-party `burn` kernels are written.

# Custom cubec<sub>l</sub> Kernel

In this section, you will learn how to create your own custom operation by writing your own kernel with the cubec<sub>l</sub> compiler frontend. We will take the example of a common workflow in the deep learning field, where we create a kernel to fuse multiple operations together. Note that `burn` does this automatically, but a manual implementation might be more efficient in some cases. We will fuse a `matmul` kernel followed by an addition and the ReLU activation function, which is commonly found in various models. All the code can be found under the [examples directory](#).

## Custom Backend Trait

First, we need to determine the type signature of our newly created operation by defining our custom backend traits. As we will use the associated type `TensorPrimitive` of the `Backend` trait, which encapsulates the underlying tensor implementation of the backend, we will use a type alias to avoid the ugly disambiguation with associated types.

```
/// We create our own Backend trait that extends the Burn backend trait.
pub trait Backend: burn::tensor::backend::Backend {
    fn fused_matmul_add_relu(
        lhs: FloatTensor<Self>,
        rhs: FloatTensor<Self>,
        bias: FloatTensor<Self>,
    ) -> FloatTensor<Self>;
}

/// We create our own AutodiffBackend trait that extends the Burn autodiff
/// backend trait.
pub trait AutodiffBackend: Backend + burn::tensor::backend::AutodiffBackend {}
```

In our project, we can use these traits instead of the `burn::tensor::backend::{Backend, AutodiffBackend}` traits provided by Burn. Burn's user APIs typically make use of the `Tensor` struct rather than dealing directly with primitive tensor types. Therefore, we can encapsulate our newly defined backend traits with functions that expose new operations while maintaining a consistent API.

```

/// We define our custom implementation using the added function on our custom
backend.
pub fn matmul_add_relu_custom<B: Backend>(
    lhs: Tensor<B, 3>,
    rhs: Tensor<B, 3>,
    bias: Tensor<B, 3>,
) -> Tensor<B, 3> {
    let output = B::fused_matmul_add_relu(
        lhs.into_primitive().tensor(),
        rhs.into_primitive().tensor(),
        bias.into_primitive().tensor(),
    );

    Tensor::from_primitive(TensorPrimitive::Float(output))
}

/// We define a reference implementation using basic tensor operations.
pub fn matmul_add_relu_reference<B: Backend>(
    lhs: Tensor<B, 3>,
    rhs: Tensor<B, 3>,
    bias: Tensor<B, 3>,
) -> Tensor<B, 3> {
    let x = lhs.matmul(rhs) + bias;

    activation::relu(x)
}

```

Note that we also provide a reference implementation for testing purposes, which allows us to easily validate our new implementation. While not mandatory, having a reference implementation can be valuable, especially in projects where creating a reference implementation solely using basic tensor operations is feasible.

## Forward Kernel

Now, let's proceed to write the fused kernel using the `cubec1` compiler frontend. To keep things simple, we'll create a straightforward matmul kernel without employing any intricate techniques. We won't delve into the details of the `cube` macro, but if you're interested to learn more, please see [cubec1 Book](#). The actual matmul, add and relu computations are found at the end, after an extensive prelude that serves to correctly map each compute unit to the data it is responsible for, with support for batches.

```

use cubec1::{cube, prelude::*};

#[cube(launch)]
pub fn fused_matmul_add_relu_kernel<F: Float>(
    lhs: &Tensor<F>,
    rhs: &Tensor<F>,
    bias: &Tensor<F>,
    output: &mut Tensor<F>,
) {
    let row = ABSOLUTE_POS_X;
    let col = ABSOLUTE_POS_Y;
    let batch = ABSOLUTE_POS_Z;

    let n_rows = output.shape(output.rank() - 2);
    let n_cols = output.shape(output.rank() - 1);
    let dim_k = rhs.shape(rhs.rank() - 1);

    if row >= n_rows || col >= n_cols {
        return;
    }

    let offset_output = batch * n_rows * n_cols;
    let mut offset_lhs = 0;
    let mut offset_rhs = 0;

    let batch_dims = output.rank() - 2;
    for dim in 0..batch_dims {
        offset_lhs += offset_output / output.stride(dim) % lhs.shape(dim) *
lhs.stride(dim);
        offset_rhs += offset_output / output.stride(dim) % rhs.shape(dim) *
rhs.stride(dim);
    }

    let mut sum = F::new(0.0);
    for k in 0..dim_k {
        let lhs_index = row * dim_k + k;
        let rhs_index = k * n_cols + col;

        sum += lhs[offset_lhs + lhs_index] * rhs[offset_rhs + rhs_index];
    }

    let out_index = row * n_cols + col;
    let index = offset_output + out_index;

    output[index] = F::max(sum + bias[index], F::new(0.0));
}

```

Now, let's move on to the next step, which involves implementing the remaining code to launch the kernel. We'll go into implementing our custom backend trait for the generic JIT backend. This automatically implements the trait for `burn-cuda`, `burn-wgpu` as well as `fusion`.

```

/// Implement our custom backend trait for the generic `JitBackend`.
impl<R: JitRuntime, F: FloatElement, I: IntElement> Backend for JitBackend<R,
F, I> {
    fn fused_matmul_add_relu(
        lhs: FloatTensor<Self>,
        rhs: FloatTensor<Self>,
        bias: FloatTensor<Self>,
    ) -> FloatTensor<Self> {
        // Define cube dim, hardcoded for simplicity.
        let cube_dim = CubeDim { x: 16, y: 16, z: 1 };

        lhs.assert_is_on_same_device(&rhs);
        lhs.assert_is_on_same_device(&bias);

        // For simplicity, make sure each tensor is continuous.
        let lhs = into_contiguous(lhs);
        let rhs = into_contiguous(rhs);
        let bias = into_contiguous(bias);

        // Get the matmul relevant shapes.
        let ndims = lhs.shape.num_dims();
        let num_rows = lhs.shape.dims[ndims - 2];
        let num_cols = rhs.shape.dims[ndims - 1];

        // Compute shape of output, while tracking number of batches.
        let mut num_batches = 1;
        let mut shape_out = vec![0; ndims];
        for i in shape_out.clone().into_iter().take(ndims - 2) {
            shape_out[i] = usize::max(lhs.shape.dims[i], rhs.shape.dims[i]);
            num_batches *= shape_out[i];
        }
        shape_out[ndims - 2] = num_rows;
        shape_out[ndims - 1] = num_cols;
        let shape_out = Shape::from(shape_out);

        // Create a buffer for the output tensor.
        let buffer = lhs
            .client
            .empty(shape_out.num_elements() * core::mem::size_of::<F>());

        // Create the output tensor primitive.
        let output =
            JitTensor::new_contiguous(lhs.client.clone(), lhs.device.clone(),
            shape_out, buffer);

        // Declare the wgsml workgroup with the number of cubes in x, y and z.
        let cubes_needed_in_x = f32::ceil(num_rows as f32 / cube_dim.x as f32)
as u32;
        let cubes_needed_in_y = f32::ceil(num_cols as f32 / cube_dim.y as f32)
as u32;
        let cube_count =
            CubeCount::Static(cubes_needed_in_x, cubes_needed_in_y, num_batches

```

```

as u32);

    // Execute lazily the kernel with the launch information and the given
    buffers. For
    // simplicity, no vectorization is performed
    fused_matmul_add_relu_kernel::launch::<F, R>(
        &lhs.client,
        cube_count,
        cube_dim,
        lhs.as_tensor_arg(1),
        rhs.as_tensor_arg(1),
        bias.as_tensor_arg(1),
        output.as_tensor_arg(1),
    );

    // Return the output tensor.
    output
}
}

```

In the preceding code block, we demonstrated how to launch the kernel that modifies the correct buffer. It's important to note that Rust's mutability safety doesn't apply here; the context has the capability to execute any mutable operation on any buffer. While this isn't a problem in the previous scenario where we only modify the newly created output buffer, it is wise to keep this in mind.

## Backward

Now that the custom backend trait is implemented for the JIT backend, you can use it to invoke the `matmul_add_relu_custom` function. However, calculating gradients is not yet possible at this stage. If your use case does not extend beyond inference, there is no need to implement any of the following code.

For the backward pass, we will leverage the backend implementation from `burn-autodiff`, which is actually generic over the backend. Instead of crafting our own `cubec1` kernel for the backward pass, we will use our fused kernel only for the forward pass, and compute the gradient using basic operations.

```

// Implement our custom backend trait for any backend that also implements our
custom backend trait.
impl<B: Backend, C: CheckpointStrategy> Backend for Autodiff<B, C> {
    fn fused_matmul_add_relu(
        lhs: FloatTensor<Self>,
        rhs: FloatTensor<Self>,
        bias: FloatTensor<Self>,
    ) -> FloatTensor<Self> {
        // Create our zero-sized type that will implement the Backward trait.
        #[derive(Debug)]
        struct FusedMatmulAddReluBackward;

        // Implement the backward trait for the given backend B, the node
gradient
        // with three other gradients to calculate (lhs, rhs, and bias).
        impl<B: Backend> Backward<B, 3> for FusedMatmulAddReluBackward {
            // Our state that we must build during the forward pass to compute
the backward pass.
            //
            // Note that we could improve the performance further by only
keeping the state of
            // tensors that are tracked, improving memory management, but for
simplicity, we avoid
            // that part.
            type State = (NodeID, NodeID, FloatTensor<B>, Shape);

            fn backward(
                self,
                ops: Ops<Self::State, 3>,
                grads: &mut Gradients,
                checkpointer: &mut Checkpointer,
            ) {
                // Get the nodes of each variable.
                let [node_lhs, node_rhs, node_bias] = ops.parents;
                // Fetch the gradient for the current node.
                let grad = grads.consume::(&ops.node);

                // Set our state.
                let (lhs_state, rhs_state, output, shape_bias) = ops.state;
                let lhs = checkpointer.retrieve_node_output(lhs_state);
                let rhs = checkpointer.retrieve_node_output(rhs_state);

                // Fetch shapes of our tensor to support broadcasting.
                let shape_lhs = B::float_shape(&lhs);
                let shape_rhs = B::float_shape(&rhs);

                // Compute the gradient of the output using the already
existing `relu_backward`
                // function in the basic Burn backend trait.
                let grad_output = B::relu_backward(output, grad);

                // Compute the lhs gradient, which is the derivative of matmul

```

```

with support for
    // broadcasting.
    let grad_lhs = broadcast_shape::<B>(
        B::float_matmul(grad_output.clone(),
            B::float_transpose(rhs)),
        &shape_lhs,
    );
    // Compute the rhs gradient, which is the derivative of matmul
with support for
    // broadcasting.
    let grad_rhs = broadcast_shape::<B>(
        B::float_matmul(B::float_transpose(lhs),
            grad_output.clone()),
        &shape_rhs,
    );
    // The add derivative is only 1, so we just need to support
broadcasting to
    // compute the bias gradient.
    let grad_bias = broadcast_shape::<B>(grad_output, &shape_bias);

    // Register the gradient for each variable based on whether
they are marked as
    // `tracked`.
    if let Some(node) = node_bias {
        grads.register::<B>(node.id, grad_bias);
    }
    if let Some(node) = node_lhs {
        grads.register::<B>(node.id, grad_lhs);
    }
    if let Some(node) = node_rhs {
        grads.register::<B>(node.id, grad_rhs);
    }
}

}

    // Prepare a stateful operation with each variable node and
corresponding graph.
    //
    // Each node can be fetched with `ops.parents` in the same order as
defined here.
    match FusedMatmulAddReluBackward
        .prepare::<C>([lhs.node.clone(), rhs.node.clone(),
bias.node.clone()])
        // Marks the operation as compute bound, meaning it will save its
        // state instead of recomputing itself during checkpointing
        .compute_bound()
        .stateful()
    {
        OpsKind::Tracked(mut prep) => {
            // When at least one node is tracked, we should register our
backward step.

            // The state consists of what will be needed for this

```



```

operation's backward pass.
    // Since we need the parents' outputs, we must checkpoint their
ids to retrieve
    // their node output at the beginning of the backward pass. We
can also save
    // utility data such as the bias shape. If we also need this
operation's output,
    // we can either save it in the state or recompute it.
    // during the backward pass. Here we choose to save it in the
state because it's a
    // compute bound operation.
    let lhs_state = prep.checkpoint(&lhs);
    let rhs_state = prep.checkpoint(&rhs);
    let bias_shape = B::float_shape(&bias.primitive);

    let output = B::fused_matmul_add_relu(
        lhs.primitive.clone(),
        rhs.primitive.clone(),
        bias.primitive,
    );

    let state = (lhs_state, rhs_state, output.clone(), bias_shape);

    prep.finish(state, output)
}
OpsKind::UnTracked(prep) => {
    // When no node is tracked, we can just compute the original
operation without
    // keeping any state.
    let output = B::fused_matmul_add_relu(lhs.primitive,
rhs.primitive, bias.primitive);
    prep.finish(output)
}
}
}
}
}

```

The previous code is self-documented to make it clearer, but here is what it does in summary:

We define `fused_matmul_add_relu` within `Autodiff<B>`, allowing any autodiff-decorated backend to benefit from our implementation. In an autodiff-decorated backend, the forward pass must still be implemented. This is achieved using a comprehensive match statement block where computation is delegated to the inner backend, while keeping track of a state. The state comprises any information relevant to the backward pass, such as input and output tensors, along with the bias shape. When an operation isn't tracked (meaning there won't be a backward pass for this specific operation in the graph), storing a state becomes unnecessary, and we simply perform the forward computation.

The backward pass uses the gradient obtained from the preceding node in the computation graph. It calculates the derivatives for `relu` (`relu_backward`), `add` (no operation is required here, as the derivative is one), and `matmul` (another `matmul` with transposed inputs). This results in gradients for both input tensors and the bias, which are registered for consumption by subsequent operation nodes.

The only remaining part is to implement our autodiff-decorated backend trait for our JIT Backend.

```
impl<R: JitRuntime, F: FloatElement, I: IntElement> AutodiffBackend
  for Autodiff<JitBackend<R, F, I>>
{
}
```

## Conclusion

In this guide, we've implemented a fused kernel using the `cubec1` compiler frontend, enabling execution on any GPU and any `cubec1` backend. By delving into the inner workings of both the JIT backend and the autodiff backend, we've gained a deeper understanding of these systems.

While extending a backend may be harder than working with straightforward tensors, the benefits can be worth it. This approach enables the crafting of custom models with greater control over execution, which can potentially greatly enhance the performance of your models.

As we conclude this guide, we hope that you have gained insights into Burn's world of backend extensions, and that it will help you to unleash the full potential of your projects.

# Custom WGPU Kernel

In this section, you will learn how to create your own custom operation by writing your own kernel with the WGPU backend. We will take the example of a common workflow in the deep learning field, where we create a kernel to fuse multiple operations together. Note that `burn` does this automatically, but a manual implementation might be more efficient in some cases. We will fuse a `matmul` kernel followed by an addition and the `ReLU` activation function, which is commonly found in various models. All the code can be found under the [examples directory](#).

## Custom Backend Trait

First, we need to determine the type signature of our newly created operation by defining our custom backend traits. As we will use the associated type `TensorPrimitive` of the `Backend` trait, which encapsulates the underlying tensor implementation of the backend, we will use a type alias to avoid the ugly disambiguation with associated types.

```
/// We create our own Backend trait that extends the Burn backend trait.
pub trait Backend: burn::tensor::backend::Backend {
    fn fused_matmul_add_relu(
        lhs: FloatTensor<Self>,
        rhs: FloatTensor<Self>,
        bias: FloatTensor<Self>,
    ) -> FloatTensor<Self>;
}

/// We create our own AutodiffBackend trait that extends the Burn autodiff
/// backend trait.
pub trait AutodiffBackend: Backend + burn::tensor::backend::AutodiffBackend {}
```

In our project, we can use these traits instead of the `burn::tensor::backend::{Backend, AutodiffBackend}` traits provided by Burn. Burn's user APIs typically make use of the `Tensor` struct rather than dealing directly with primitive tensor types. Therefore, we can encapsulate our newly defined backend traits with functions that expose new operations while maintaining a consistent API.

```

/// We define our custom implementation using the added function on our custom
backend.
pub fn matmul_add_relu_custom<B: Backend>(
    lhs: Tensor<B, 3>,
    rhs: Tensor<B, 3>,
    bias: Tensor<B, 3>,
) -> Tensor<B, 3> {
    let output = B::fused_matmul_add_relu(
        lhs.into_primitive().tensor(),
        rhs.into_primitive().tensor(),
        bias.into_primitive().tensor(),
    );

    Tensor::from_primitive(TensorPrimitive::Float(output))
}

/// We define a reference implementation using basic tensor operations.
pub fn matmul_add_relu_reference<B: Backend>(
    lhs: Tensor<B, 3>,
    rhs: Tensor<B, 3>,
    bias: Tensor<B, 3>,
) -> Tensor<B, 3> {
    let x = lhs.matmul(rhs) + bias;

    activation::relu(x)
}

```

Note that we also provide a reference implementation for testing purposes, which allows us to easily validate our new implementation. While not mandatory, having a reference implementation can be valuable, especially in projects where creating a reference implementation solely using basic tensor operations is feasible.

## Forward Kernel

Now, let's proceed to write the fused kernel using the WGSL shading language. To keep things simple, we'll create a straightforward matmul kernel without employing any intricate techniques. Although we won't delve into the details of the WGSL syntax, as it falls beyond the scope of this guide, we still provide the implementation below for readers who are curious. The actual matmul, add and relu computations are found at the end, after an extensive overhead whose use is to correctly map each compute unit to the data it is responsible of, with support for batches.

```

@group(0)
@binding(0)
var<storage, read_write> lhs: array<{{ elem }}>;

@group(0)
@binding(1)
var<storage, read_write> rhs: array<{{ elem }}>;

@group(0)
@binding(2)
var<storage, read_write> bias: array<{{ elem }}>;

@group(0)
@binding(3)
var<storage, read_write> output: array<{{ elem }}>;

@group(0)
@binding(4)
var<storage, read_write> info: array<u32>;

const BLOCK_SIZE = {{ workgroup_size_x }}u;

@compute
@workgroup_size({{ workgroup_size_x }}, {{ workgroup_size_y }}, 1)
fn main(
    @builtin(global_invocation_id) global_id: vec3<u32>,
    @builtin(local_invocation_index) local_idx: u32,
    @builtin(workgroup_id) workgroup_id: vec3<u32>,
) {
    // Indices
    let row = workgroup_id.x * BLOCK_SIZE + (local_idx / BLOCK_SIZE);
    let col = workgroup_id.y * BLOCK_SIZE + (local_idx % BLOCK_SIZE);
    let batch = global_id.z;

    // Basic information
    let dim = info[0];
    let n_rows = info[6u * dim - 1u];
    let n_cols = info[6u * dim];
    let K = info[5u * dim - 1u];

    // Returns if outside the output dimension
    if row >= n_rows || col >= n_cols {
        return;
    }

    // Calculate the corresponding offsets with support for broadcasting.
    let offset_output = batch * n_rows * n_cols;
    var offset_lhs: u32 = 0u;
    var offset_rhs: u32 = 0u;

    let batch_dims = dim - 2u;
    for (var b: u32 = 1u; b <= batch_dims; b++) {

```

```

    let stride_lhs = info[b];
    let stride_rhs = info[b + dim];
    let stride_output = info[b + 2u * dim];
    let shape_lhs = info[b + 3u * dim];
    let shape_rhs = info[b + 4u * dim];

    offset_lhs += offset_output / stride_output % shape_lhs * stride_lhs;
    offset_rhs += offset_output / stride_output % shape_rhs * stride_rhs;
}

// Basic matmul implementation
var sum = 0.0;
for (var k: u32 = 0u; k < K; k++) {
    let lhs_index = row * K + k;
    let rhs_index = k * n_cols + col;

    sum += lhs[offset_lhs + lhs_index] * rhs[offset_rhs + rhs_index];
}

let output_index = row * n_cols + col;
let index = offset_output + output_index;

// Add and ReLU
output[index] = max(sum + bias[index], 0.0);
}

```

Now, let's move on to the next step, which involves implementing the remaining code to launch the kernel. The initial part entails loading the template and populating it with the appropriate variables. The `register(name, value)` method simply replaces occurrences of `{{ name }}` in the above WGSL code with some other string before it is compiled. In order to use templating utilities, you will have to activate the `template` feature of Burn in your `cargo.toml`.

```

// Source the kernel written in WGSL.
kernel_wgsl!(FusedMatmulAddReluRaw, "./kernel.wgsl");

// Define our kernel type with cube information.
#[derive(new, Debug)]
struct FusedMatmulAddRelu<E: FloatElement> {
    cube_dim: CubeDim,
    _elem: PhantomData<E>,
}

// Implement the dynamic kernel trait for our kernel type.
impl<E: FloatElement> KernelSource for FusedMatmulAddRelu<E> {
    fn source(&self) -> SourceTemplate {
        // Extend our raw kernel with cube size information using the
        // `SourceTemplate` trait.
        FusedMatmulAddReluRaw::new()
            .source()
            .register("workgroup_size_x", self.cube_dim.x.to_string())
            .register("workgroup_size_y", self.cube_dim.y.to_string())
            .register("elem", E::type_name())
            .register("int", "i32")
    }

    fn id(&self) -> cubec::KernelId {
        cubec::KernelId::new:::<Self>().info(self.cube_dim)
    }
}

```

Subsequently, we'll go into implementing our custom backend trait for the WGPU backend. Note that we won't go into supporting the `fusion` feature flag in this tutorial, so we implement the trait for the raw `WgpuBackend` type.

```

/// Implement our custom backend trait for the existing backend `WgpuBackend`.
impl<F: FloatElement, I: IntElement> Backend for JitBackend<WgpuRuntime, F, I>
{
    fn fused_matmul_add_relu(
        lhs: FloatTensor<Self>,
        rhs: FloatTensor<Self>,
        bias: FloatTensor<Self>,
    ) -> FloatTensor<Self> {
        // Define cube dim, hardcoded for simplicity.
        let cube_dim = CubeDim { x: 16, y: 16, z: 1 };

        lhs.assert_is_on_same_device(&rhs);
        lhs.assert_is_on_same_device(&bias);

        // For simplicity, make sure each tensor is continuous.
        let lhs = into_contiguous(lhs);
        let rhs = into_contiguous(rhs);
        let bias = into_contiguous(bias);

        // Get the matmul relevant shapes.
        let ndims = lhs.shape.num_dims();
        let num_rows = lhs.shape.dims[ndims - 2];
        let num_cols = rhs.shape.dims[ndims - 1];

        // Compute shape of output, while tracking number of batches.
        let mut num_batches = 1;
        let mut shape_out = vec![0; ndims];
        for i in shape_out.clone().into_iter().take(ndims - 2) {
            shape_out[i] = usize::max(lhs.shape.dims[i], rhs.shape.dims[i]);
            num_batches *= shape_out[i];
        }
        shape_out[ndims - 2] = num_rows;
        shape_out[ndims - 1] = num_cols;
        let shape_out = Shape::from(shape_out);

        // Create a buffer for the output tensor.
        let buffer = lhs
            .client
            .empty(shape_out.num_elements() * core::mem::size_of::<F>());

        // Create the output tensor primitive.
        let output =
            JitTensor::new_contiguous(lhs.client.clone(), lhs.device.clone(),
            shape_out, buffer);

        // Create the kernel.
        let kernel = FusedMatmulAddRelu::<F>::new(cube_dim);

        // Build info buffer with tensor information needed by the kernel, such
        as shapes and strides.
        let info = build_info(&[&lhs, &rhs, &output]);
        let info_handle = lhs.client.create(bytemuck::cast_slice(&info));
    }
}

```



```

        // Declare the wgsl workgroup with the number of cubes in x, y and z.
        let cubes_needed_in_x = f32::ceil(num_rows as f32 / cube_dim.x as f32)
as u32;
        let cubes_needed_in_y = f32::ceil(num_cols as f32 / cube_dim.y as f32)
as u32;
        let cube_count =
            CubeCount::Static(cubes_needed_in_x, cubes_needed_in_y, num_batches
as u32);

        // Execute lazily the kernel with the launch information and the given
        buffers.
        lhs.client.execute(
            Box::new(SourceKernel::new(kernel, cube_dim)),
            cube_count,
            vec![
                lhs.handle.binding(),
                rhs.handle.binding(),
                bias.handle.binding(),
                output.handle.clone().binding(),
                info_handle.binding(),
            ],
        );

        // Return the output tensor.
        output
    }
}

```

In the preceding code block, we demonstrated how to launch the kernel that modifies the correct buffer. It's important to note that Rust's mutability safety doesn't apply here; the context has the capability to execute any mutable operation on any buffer. While this isn't a problem in the previous scenario where we only modify the newly created output buffer, it is wise to keep this in mind.

## Backward

Now that the custom backend trait is implemented for the WGPU backend, you can use it to invoke the `matmul_add_relu_custom` function. However, calculating gradients is not yet possible at this stage. If your use case does not extend beyond inference, there is no need to implement any of the following code.

For the backward pass, we will leverage the backend implementation from `burn-autodiff`, which is actually generic over the backend. Instead of crafting our own WGSL kernel for the backward pass, we will use our fused kernel only for the forward pass, and compute the

gradient using basic operations.

```

// Implement our custom backend trait for any backend that also implements our
custom backend trait.
//
// Note that we could implement the backend trait only for the Wgpu backend
instead of any backend that
// also implements our own API. This would allow us to call any function only
implemented for Wgpu
// and potentially call a custom kernel crafted only for this task.
impl<B: Backend, C: CheckpointStrategy> Backend for Autodiff<B, C> {
    fn fused_matmul_add_relu(
        lhs: FloatTensor<Self>,
        rhs: FloatTensor<Self>,
        bias: FloatTensor<Self>,
    ) -> FloatTensor<Self> {
        // Create our zero-sized type that will implement the Backward trait.
        #[derive(Debug)]
        struct FusedMatmulAddReluBackward;

        // Implement the backward trait for the given backend B, the node
gradient
        // with three other gradients to calculate (lhs, rhs, and bias).
        impl<B: Backend> Backward<B, 3> for FusedMatmulAddReluBackward {
            // Our state that we must build during the forward pass to compute
the backward pass.
            //
            // Note that we could improve the performance further by only
keeping the state of
            // tensors that are tracked, improving memory management, but for
simplicity, we avoid
            // that part.
            type State = (NodeID, NodeID, FloatTensor<B>, Shape);

            fn backward(
                self,
                ops: Ops<Self::State, 3>,
                grads: &mut Gradients,
                checkpointer: &mut Checkpointer,
            ) {
                // Get the nodes of each variable.
                let [node_lhs, node_rhs, node_bias] = ops.parents;
                // Fetch the gradient for the current node.
                let grad = grads.consume::<B>(&ops.node);

                // Set our state.
                let (lhs_state, rhs_state, output, shape_bias) = ops.state;
                let lhs = checkpointer.retrieve_node_output(lhs_state);
                let rhs = checkpointer.retrieve_node_output(rhs_state);

                // Fetch shapes of our tensor to support broadcasting.
                let shape_lhs = B::float_shape(&lhs);
                let shape_rhs = B::float_shape(&rhs);

```

```

        // Compute the gradient of the output using the already
existing `relu_backward`
        // function in the basic Burn backend trait.
        let grad_output = B::relu_backward(output, grad);

        // Compute the lhs gradient, which is the derivative of matmul
with support for
        // broadcasting.
        let grad_lhs = broadcast_shape::<B>(
            B::float_matmul(grad_output.clone(),
B::float_transpose(rhs)),
            &shape_lhs,
        );
        // Compute the rhs gradient, which is the derivative of matmul
with support for
        // broadcasting.
        let grad_rhs = broadcast_shape::<B>(
            B::float_matmul(B::float_transpose(lhs),
grad_output.clone()),
            &shape_rhs,
        );
        // The add derivative is only 1, so we just need to support
broadcasting to
        // compute the bias gradient.
        let grad_bias = broadcast_shape::<B>(grad_output, &shape_bias);

        // Register the gradient for each variable based on whether
they are marked as
        // `tracked`.
        if let Some(node) = node_bias {
            grads.register::<B>(node.id, grad_bias);
        }
        if let Some(node) = node_lhs {
            grads.register::<B>(node.id, grad_lhs);
        }
        if let Some(node) = node_rhs {
            grads.register::<B>(node.id, grad_rhs);
        }
    }
}

    // Prepare a stateful operation with each variable node and
corresponding graph.
    //
    // Each node can be fetched with `ops.parents` in the same order as
defined here.
    match FusedMatmulAddReluBackward
        .prepare::<C>([lhs.node.clone(), rhs.node.clone(),
bias.node.clone()])
    {
        // Marks the operation as compute bound, meaning it will save its
// state instead of recomputing itself during checkpointing
        .compute_bound()
        .stateful()
    }
}

```

```

{
  OpsKind::Tracked(mut prep) => {
    // When at least one node is tracked, we should register our
backward step.

    // The state consists of what will be needed for this
operation's backward pass.
    // Since we need the parents' outputs, we must checkpoint their
ids to retrieve their node
    // output at the beginning of the backward. We can also save
utility data such as the bias shape
    // If we also need this operation's output, we can either save
it in the state or recompute it
    // during the backward pass. Here we choose to save it in the
state because it's a compute bound operation.
    let lhs_state = prep.checkpoint(&lhs);
    let rhs_state = prep.checkpoint(&rhs);
    let bias_shape = B::float_shape(&bias.primitive);

    let output = B::fused_matmul_add_relu(
      lhs.primitive.clone(),
      rhs.primitive.clone(),
      bias.primitive,
    );

    let state = (lhs_state, rhs_state, output.clone(), bias_shape);

    prep.finish(state, output)
  }
  OpsKind::UnTracked(prepare) => {
    // When no node is tracked, we can just compute the original
operation without
    // keeping any state.
    let output = B::fused_matmul_add_relu(lhs.primitive,
rhs.primitive, bias.primitive);
    prep.finish(output)
  }
}
}
}

```

The previous code is self-documented to make it clearer, but here is what it does in summary.

We define `fused_matmul_add_relu` within `Autodiff<B>`, allowing any autodiff-decorated backend to benefit from our implementation. In an autodiff-decorated backend, the forward pass must still be implemented. This is achieved using a comprehensive match statement block where computation is delegated to the inner backend, while keeping track of a state. The state comprises any information relevant to the backward pass, such as input and output tensors, along with the bias shape. When an operation isn't tracked (meaning there

won't be a backward pass for this specific operation in the graph), storing a state becomes unnecessary, and we simply perform the forward computation.

The backward pass uses the gradient obtained from the preceding node in the computation graph. It calculates the derivatives for `relu` (`relu_backward`), `add` (no operation is required here, as the derivative is one), and `matmul` (another `matmul` with transposed inputs). This results in gradients for both input tensors and the bias, which are registered for consumption by subsequent operation nodes.

The only remaining part is to implement our autodiff-decorated backend trait for our WGPU Backend.

```
impl<G: GraphicsApi, F: FloatElement, I: IntElement> AutodiffBackend for
Autodiff<WgpuBackend<G, F, I>>
{
}
```

## Conclusion

In this guide, we've implemented a fused kernel using the WGPU backend, enabling execution on any GPU. By delving into the inner workings of both the WGPU backend and the autodiff backend, we've gained a deeper understanding of these systems.

While extending a backend may be harder than working with straightforward tensors, the benefits can be worth it. This approach enables the crafting of custom models with greater control over execution, which can potentially greatly enhance the performance of your models.

As we conclude this guide, we hope that you have gained insights into Burn's world of backend extensions, and that it will help you to unleash the full potential of your projects.

# No Standard Library

In this section, you will learn how to run an onnx inference model on an embedded system, with no standard library support on a Raspberry Pi Pico. This should be universally applicable to other platforms. All the code can be found under the [examples directory](#).

## Step-by-Step Guide

Let's walk through the process of running an embedded ONNX model:

### Setup

Follow the [embassy guide](#) for your specific environment. Once setup, you should have something similar to the following.

```
./inference
├── Cargo.lock
├── Cargo.toml
├── build.rs
├── memory.x
├── src
│   └── main.rs
```

Some other dependencies have to be added

```
[dependencies]
embedded-alloc = "0.5.1" # Only if there is no default allocator for your chip
burn = { version = "0.15", default-features = false, features = ["ndarray"] } #
Backend must be ndarray
```

```
[build-dependencies]
burn-import = { version = "0.14" } # Used to auto generate the rust code to
import the model
```

### Import the Model

Follow the directions to [import models](#).

Use the following ModelGen config

```
ModelGen::new()
  .input(my_model)
  .out_dir("model/")
  .record_type(RecordType::Bincode)
  .embed_states(true)
  .run_from_script();
```

## Global Allocator

First define a global allocator (if you are on a no\_std system without alloc).

```
use embedded_alloc::Heap;

#[global_allocator]
static HEAP: Heap = Heap::empty();

#[embassy_executor::main]
async fn main(_spawner: Spawner) {
    {
        use core::mem::MaybeUninit;
        const HEAP_SIZE: usize = 100 * 1024; // This is dependent on the model
        size in memory.
        static mut HEAP_MEM: [MaybeUninit<u8>; HEAP_SIZE] =
        [MaybeUninit::uninit(); HEAP_SIZE];
        unsafe { HEAP.init(HEAP_MEM.as_ptr() as usize, HEAP_SIZE) }
    }
}
```

## Define Backend

We are using ndarray, so we just need to define the NdArray backend as usual

```
use burn::{backend::NdArray, tensor::Tensor};

type Backend = NdArray<f32>;
type BackendDeice = <Backend as burn::tensor::backend::Backend>::Device;
```

Then inside the main function add



```
use your_model::Model;

// Get a default device for the backend
let device = BackendDeice::default();

// Create a new model and load the state
let model: Model<Backend> = Model::default();
```

## Running the Model

To run the model, just call it as you would normally

```
// Define the tensor
let input = Tensor::<Backend, 2>::from_floats([[input]], &device);

// Run the model on the input
let output = model.forward(input);
```

## Conclusion

Running a model in a no\_std environment is pretty much identical to a normal environment. All that is needed is a global allocator.