

# JAX – Now and Future

Nasy

Jul 16, 2023

## Contents

### 1 Introduction

#### 1.1 What is JAX?

##### 1.1.1 Acronym

**Originally** Just After eXecution

**Now** JAX is Autograd (automatic obtaining of the gradient function through differentiation of a function) an XLA (accelerated linear algebra).

**Fact** JAX is JAX

ref: <https://github.com/google/jax/discussions/9019>

##### 1.1.2 Design

- Follow `numpy` as closely as possible
- Works with various existing frameworks (PyTorch, Tensorflow)
- Immutable and purely functional
- Asynchronous dispatch
- Core: `jit`, `vmap`, `grad`, `pmap`

## 1.2 Who use JAX?

- Lab
  - Google Brain & Deepmind
  - Google Research
- Models
  - ViT
  - Big Vision
  - AlphaFold
  - MLP-Mixer
  - T5X
  - PaLM
  - ...
- Awesome-jax: <https://github.com/n2cholas/awesome-jax>

## 1.3 Basic Usage

### 1.3.1 JAX and NumPy

1. JAX is accelerated NumPy

```
[[python import jax import jax.numpy as jnp import numpy as np
from rich import print

print("numpy:", np.asarray([1, 2, 3])) print("jax:", jnp.asarray([1, 2,
3]))

[[python print("jax->numpy:", np.std(jnp.arange(10))) print("numpy-
>jax:", jnp.std(np.arange(10)))
```

2. Difference

JAX is designed to be functional, as in functional programming.

```
[[python x = np.arange(10) jx = jnp.arange(10) print(f"Original: x=")
x[0] = 1 print(f"Inplace replace: x=")
cannot jx[0] = 1 jx = jx.at[0].set(1)
```

```
def inplace_set(x, i, v) : x[i] = v return x
```

```
inplace_set(x, 2, 10)
```

```
print(f"Inplace replace 2: x=")
```

### 1.3.2 JIT

Using a just-in-time (JIT) compilation decorator, sequences of operations can be optimized together and run at once.

ref: [https://jax.readthedocs.io/en/latest/\\_autosummary/jax.jit.html#jax.jit](https://jax.readthedocs.io/en/latest/_autosummary/jax.jit.html#jax.jit)  
ref2, tutorial: <https://jax.readthedocs.io/en/latest/jax-101/02-jitting.html>

Here is an example to JIT a function.

```
[]python def func(x: Array) -> Array: ...
```

```
jit_func = jax.jit(func)
```

or

```
@jax.jit def func(x: Array) -> Array: ...
```

In most of time, you only need add the jit function in the outer most function.

```
[]python import jax import jax.numpy as jnp import rich import time
def model(params, x): return params["w"] * x + params["b"]
def loss_function(params, x, y) : return ((model(params, x)-y)**2).mean()
def train_step(params, x, y, lr) : grads = jax.grad(loss_function)(params, x, y) return "w" : params[
def train(x, y, lr, num_steps) : params = "w" : 0.0, "b" : 0.0 for i in range(num_steps) :
params = train_step(params, x, y, lr) return params
train_step_jit = jax.jit(train_step)
def train_jit(x, y, lr, num_steps) : params = "w" : 0.0, "b" : 0.0 for i in range(num_steps) :
params = train_step_jit(params, x, y, lr) return params
xs = jnp.array([1, 2, 3, 4, 5]) / 2 ys = jnp.array([2, 4, 6, 8, 10]) / 2
t1 = time.time() p1 = train(xs, ys, 0.01, 50) t2 = time.time()
t3 = time.time() p2 = train_jit(xs, ys, 0.01, 50) t4 = time.time()
rich.print("no jit:", p1, t2 - t1) rich.print("jit:", p2, t4 - t3)
```

### 1.3.3 autograd

JAX use `jax.grad` and `jax.value_and_grad` to get the gradient of a function.

```
[]python def f(x): return x ** 3
print(jax.grad(f)(10.0)) 3x2 print(jax.value_and_grad(f)(10.0))(x3, 3x2) print(jax.grad(jax.grad(f))
6x
```

You can also do partial differentiation with some structure data in JAX.

```
[]python def f(xy): return xy["x"] ** xy["y"]
```

```
x = jnp.array(2.) y = jnp.array(3.)
```

```
print(jax.value_and_grad(f)("x" : x, "y" : y)) Higher order derivatives print(jax.hessian(f)("x" : x,
```

### 1.3.4 vmap and pmap

JAX for single-program, multiple-data (SPMD).

`jax.vmap(f)(x)`, where the shape of `x` is `batch_size, ...`

Here is an example for vmap

In python we want to calculate the gradient for x and y, however, our x and y is batched.
 

```
@jax.grad def f2(xy): x, y = xy return x ** y
```

```
xs = jnp.array([2., 2.]) ys = jnp.array([3., 3.])
```

```
for-loop grads = []
for x, y in zip(xs, ys):
    grads.append(f2((x, y)))
```

```
print("For-loop:", grads)
```

```
vmap vmap_grads = jax.vmap(f2) print("vmap:", vmap_grads((xs, ys)))
```

How about pmap?

`pmap` is like `vmap`, but parallel evaluate the function on different devices.

`jax.pmap(f)(x)`, where the shape of `x` is `devices, ...`

And you can use both `pmap` and `vmap`, like `jax.pmap(jax.vmap(f))` the shape will be `devices, batch_size, ...`

### 1.3.5 Performance

Here is a comparison between numpy, jax and pytorch.

```
python import time import torch np.random.seed(42)
```

```
arr = np.random.random(1000000).reshape(-1, 1000).astype("float32")
```

```
* 10 jrr = jnp.array(arr)
```

```
def func(x): return x @ x @ x @ x @ x @ x @ x @ x @ x @ x
```

```
def func(x): return (f $unc(x) *_{\mathcal{F}}$   $unc(x) +_{\mathcal{F}}$   $unc(x) *_{\mathcal{F}}$   $unc(x)$ )/f $unc(x)$ 
```

```
import time
```

```
numpy time t1 = time.time() func(arr) t2 = time.time()
```

```
jax time t3 = time.time() func(jrr).blockuntil ready() t4 = time.time()
```

```
jax.jit compile jit_func = jax.jit(func)jit_func(jrr).block_until_ready()
```

```
jax jit time t5 = time.time() jit_func(jrr).block_until_ready() t6 = time.time()
```

```
trr = torch.from_numpy(arr)
```

```
torch time func(trr) t7 = time.time() func(trr) t8 = time.time()
```

```
print("Numpy time: ", t2 - t1) print("Jax time: ", t4 - t3) print("Jax jit
```

```
time: ", t6 - t5) print("Torch time: ", t8 - t7)
```

## 2 JAX vs PyTorch in a Pipeline

## 2.1 Keras

Those who have previously used TensorFlow should be quite familiar with Keras. Currently, Keras has reached version 3.0 and includes the 'keras-core'

library, which allows for very easy switching between JAX, TensorFlow, and PyTorch.

```

[]sh run with jax KERASBACKEND = jaxpythontrain.py
run with torch KERASBACKEND = torchpythontrain.py
run with tensorflow KERASBACKEND = tensorflowpythontrain.py

```

## 2.2 Install and run

Install

```

[]sh jax pip install --upgrade "jaxlib[cuda12_pip]" --fhttps :
//storage.googleapis.com/jax-releases/jax_cuda_releases.htmljaxnnlibrarypipinstallflax
torch pip install torch

```

Run. As far as I know, both JAX and PyTorch now ship with nvidia-cuda-toolkits, so you do not need to setup LD\_LIBRARY\_PATH anymore.

```

[]sh jax python train.py
torch python train.py

```

## 2.3 Load data

JAX does not have built-in data loading utilities, so we can use both `tensorflow` or `torch` dataloader to load the dataset.

Here is a example for `torch` dataloader.

```

[]python import jax.numpy as jnp from jax.tree_util import tree_map from torch.utils import data
def collate_fn(x): return tree_map(jnp.asarray, data.default_collate(batch))
data_generator = Dataloader(dataset, collate_fn = collate_fn, batch_size =
128, shuffle = False, num_workers = 2,)

```

And you can find others in the `jax` document. [https://jax.readthedocs.io/en/latest/advanced\\_guide.html](https://jax.readthedocs.io/en/latest/advanced_guide.html)

## 2.4 Define and initialize model

### 2.4.1 PyTorch

```

[]python import torch from torch import nn
class TM(nn.Module): """Torch Model."""
def __init__(self, in=100, h1=300, h2=200, h3=100): super().__init__() self.l1=nn.Linear(in, h1) self.bn1=nn.BatchNorm1d(h1) self.dp1=nn.Dropout
def forward(self, x): x = torch.relu(self.l1(x)) x = self.bn1(x) x = self.dp1(x)
x = torch.relu(self.l2(x)) x = torch.relu(self.l3(x)) return self.out(x)
tm = TM() print(tm) rich.print(jax.tree_map(lambda x : x.shape, dict(tm.state_dict())))

```

### 2.4.2 JAX w/ flax

```
python import jax import jax.numpy as jnp from flax import linen
class JM(linen.Module): """JAX and flax model."""
    h1 = 300 h2 = 200 h3 = 100
    @linen.compact def call(self, x, training=True): x = linen.Dense(self.h1)(x) x = linen.BatchNorm()(x, use_running_average=not training)
    jm = JM() variables = jm.init(jax.random.key(42), random key jnp.ones((1, 100)), input (Batch, Features) training=False, training mode )
    rich.print(jax.tree_map(lambda x : x.shape, variables)) print(jm.tabulate(jax.random.key(42), jnp.ones((1, 100)), input (Batch, Features) training=False, compute_flops = True, compute_vjp_flops = True))
```

## 2.5 TrainState and Train loop

Usually, we need to store the state of the model, like the parameters, the optimizer, the learning rate scheduler, etc.

### 2.5.1 PyTorch

```
python import torch import torch.nn as nn from typing import NamedTuple
def train_loop(conf : dict, dataloader) : """The pytorch training loop.""" model = TM(conf)
    loss_fn = nn.CrossEntropyLoss()
    optim = torch.optim.Adam(model.parameters(), lr=conf["lr"])
    state = TrainState( loss=0, state=model.state_dict(), optim_state = optim.state_dict(), step = 0, metric = 0, )
    model.to("cuda")
    def train_step(batch, ys) : batch, ys = batch.to("cuda"), ys.to("cuda")
    forward loss = loss_fn(model(batch), ys) gradloss.backward() update params optim.step() optim.zero_grad()
    return loss
    for e in range(conf["epochs"]): model.train() for i, (batch, ys) in dataloader: loss = train_step(batch, ys)
    eval model.eval() metric = ... compute metric
    state = state.replace(loss = loss.item(), state = model.state_dict(), optim_state = optim.state_dict(), step = state.step + 1, metric = metric, )
    save torch.save(state, "model.pt")
```

### 2.5.2 JAX

```
python import jax import jax.numpy as jnp import optax from flax.core.scope
import Collection from flax.training import train_state import pickle
```

```

class TrainState(train_state.TrainState) : """Training states."""
    default_apply_fn = the_model_forward_function
    param_store = optim_step_training_step
    our_custom_batch_stats : Collection
    our_metrics_loss : jax.Array
    metric : jax.Array

    def create_train_state(conf : dict) : """Create initial training state."""
    model = JM(conf)
    @jax.jit def init() -> Collection: return model.init(jax.random.key(42),
    jnp.ones((1, 100)), training=False)
    variables = init()
    return TrainState.create(apply_fn = model.apply, params = variables["params"], tx =
    optax.adamw(conf["lr"]), model_state = Collection(), loss = jnp.inf, metric =
    0.0, batch_stats = variables["batch_stats"], step = 0, )
    @jax.jit def train_step(state, rng, batch, ys) :
    @jax.jit optional since we already have jit outside train_step
    @jax.value_and_grad def loss_fn(params) :
    return optax.cross_entropy_loss(state.apply_fn("params" : state.params, "batch_stats" : state.batch_stats,
    True, rngs = "dropout" : rng, mutable = ["batch_stats"], ), ys, ).mean()
    loss, grad = loss_fn(state.params)
    the_step, params, tx, states will automatically be updated
    return state.replace(loss = loss, grad = grad, loss = loss)
    @jax.jit def eval_step(state, batch, ys) :
    value = state.apply_fn("params" : state.params, "batch_stats" : state.batch_stats, batch, training =
    False, )
    metric = ...return metric
    def train_loop(conf : dict, dataloader) : """Training loop."""
    state = create_train_state(conf)
    for epoch in range(conf["epochs"]):
    for batch, ys in dataloader:
    state = train_step(state, batch)
    eval_metric = eval_step(state, batch, ys)
    state = state.replace(metric = eval_metric, metric = eval_metric)
    save with open("model.pkl", "wb") as f:
    pickle.dump("params": state.params, "tx": state.tx, "batch_stats" : state.batch_stats, "metric" : state.metric, f)

```

## 3 Parallel and Distributed Computing

### 3.1 Resource

- Flax tutorial: [https://flax.readthedocs.io/en/latest/guides/parallel\\_training/index.html](https://flax.readthedocs.io/en/latest/guides/parallel_training/index.html)
- JAX tutorial: [https://jax.readthedocs.io/en/latest/notebooks/Distributed\\_arrays\\_and\\_automatic\\_parallelization.html](https://jax.readthedocs.io/en/latest/notebooks/Distributed_arrays_and_automatic_parallelization.html)

### 3.2 Transfer data between devices

```
python import jax import jax.numpy as jnp import os
os.environ["XLA_FLAGS"] = " --xla_force_host_platform_device_count = 8"

x = jnp.ones((8, 8))
1. check devices print("global devices:", jax.devices()) print("local de-
vices:", jax.local_devices()) you can specify the device type print("cpu devices :
", jax.devices("cpu"))
2. check the device of x print("x devices:", x.devices())
Put x to a specific device y = jax.device_put(x, jax.devices("cpu")[1]) Get y back to host (numpy array) z =
jax.device_get(y) print("x devices : ", x.devices()) print("y devices : ", y.devices()) print(" z type :
", type(z))
```

### 3.3 Use pmap to train with data parallel.

```
python import jax from functools import partial
@partial(jax.pmap, static_broadcasted_argnums = (0, )) def create_train_state(conf :
dict) : """ Create initial training state. """ ...
@jax.pmap def train_step(state, rng, batch, ys) : ...
@jax.jit def eval_step(state, batch, ys) :
preds = state.apply_fn((batch, 1) "params" : state.params, "batch_stats" : state.batch_stats, batch, tra
False, ) logits = jax.nn.sigmoid(preds) acc = jnp.mean((logits > 0.5) ==
ys) return acc
@partial(jax.pmap, axis_name = "batch") def eval_step_pmap(state, batch, ys) :
preds = state.apply_fn((devices, batch, 1) "params" : state.params, "batch_stats" : state.batch_stats, b
False, ) logits = jax.nn.sigmoid(preds) pacc = jax.lax.pmean((logits >
0.5) == ys, axis_name = "batch") acc = jnp.mean(pacc) return acc
```

### 3.4 Use jit with sharding.

sharding, split a large array into smaller pieces, and each piece is stored on a different device.

#### 3.4.1 Basic

```
python import jax import jax.numpy as jnp from jax.experimental import
mesh_utils from jax.sharding import PositionalSharding, NamedSharding import numpy as np
import os
os.environ["XLA_FLAGS"] = " --xla_force_host_platform_device_count = 8"
```



```

dc = jax.device_count()print(f"wehavedcdevices.")
If you want to track gpu usage install go pip install jax-smi from
jax_smiimportinitialise_trackinginitialise_tracking()
batch feature (batch, feature), here batch = Nxdevice_countxs = jax.random.normal(jax.random.k
dc * 64, 8 * dc * 64))
dmesh dmesh = mesh_utils.create_device_mesh((dc,))ordmesh = np.asarray(jax.devices()).reshape(
sharding = PositionalSharding(dmesh)
since xs shap is (batch, feature) sharding = sharding.reshape((-1, 1))
put jax across devices print(xs.devices()) xs = jax.device_put(xs, sharding)print(xs.devices())
jax.debug.visualize_arraysharding(xs)

```

### 3.4.2 Different shardings

```

[]python xs = jax.device_put(xs, sharding.reshape(1, -1))
jax.debug.visualize_arraysharding(xs)
[]python xs = jax.device_put(xs, sharding.reshape(dc//2, -1))jax.debug.visualize_arraysharding(xs)

```

### 3.4.3 Performance

```

[]python xs = jax.device_put(xs, sharding.reshape(dc//2, -1))xsc0 = jax.device_put(xs, jax.devices()[0])
cos = jax.jit(jnp.cos) cos(xsc0).block_until_ready()cos(xs).block_until_ready()
t1 = time.time() for i in range(20) : r = cos(xsc0).block_until_ready()t2 =
time.time()
t3 = time.time() for i in range(20) : rr = cos(xs).block_until_ready()t4 =
time.time()
cos = jax.pmap(jnp.cos)
pxs = xsc0.reshape(dc, -1, xsc0.shape[-1])
t5 = time.time() for i in range(20) : rrr = cos(pxs).block_until_ready()t6 =
time.time()
print("In single device: ", t2 - t1) print("In multi device: ", t4 - t3)
print("In multi device pmap: ", t6 - t5)

```

### 3.4.4 Use sharding with jit

```

[]python xs = jax.device_put(xs, sharding.reshape(dc, -1))
def f(x): return jnp.sin(x)
in with None, out with (dc, 1)
print("before:") jax.debug.visualize_arraysharding(xs)
out = jax.jit(f, in_shardings = sharding.reshape(dc, -1), out_shardings =
sharding.reshape(-1, dc))(xs)print("after : ")jax.debug.visualize_arraysharding(out)

```

## 4 JAX ecosystem

- JAX
- NN library
  - Flax (Google)
  - Haiku (Deepmind)
  - Trax (Google brain)
  - HuggingFace (Flax)
  - keras
  - jraph (GNN)
  - RLax (Deepmind, RL)
  - Coax (Microsoft, RL)
- Optimizer
  - jaxopt
  - optax
- Others
  - orbax-checkpoint
  - jax-md (molecular dynamics)
  - mpi4jax