

Goals of this talk

- 1) What features are we looking for in an ML/DL framework?
- 2) What is JAX?
- 3) What sets JAX apart from other frameworks?
- 4) How can you train Neural Networks in JAX with Flax?
- 5) Where can I continue my learning journey into JAX with Flax?





Successes of Machine and Deep Learning



AlphaFold







Explain quantum computing in simple terms

C

Quantum computing is a type of computing that uses quantum bits, or qubits, instead of ムマ classical bits to store and process information. Unlike classical bits, which can only be in one of two states (0 or 1), qubits can be in a superposition of both states at the same time. This allows quantum computers to perform certain types of calculations much more efficiently than classical computers.



fairytale book.

Sprouts in the shape of text 'Imagen' coming out of a A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.

ChatGPT

Imagen

Secret sauce: Enormous compute power

How can we make efficient use of the compute power? => DL frameworks



Basic Requirements on a Deep Learning Framework



Automatic Differentiation

 X_1

X₂

Gradients



Why yet another framework?

Frameworks follow different philosophies



Greater user experience

Best of both worlds?





Faster / Better to optimize



Properties	Ore PyTorch	
Design	Object-Oriented	Functional Programming
Speed		
Community / Doc.		
Debugging		
Code Minimalism		
Flexibility		









Properties	O PyTorch		
Design	Object-Oriented	Functiona	l Programmi
Speed	Ŧ	G	
Community / Doc.	Models	PvTorch	JAX
Debugging	GoogleNet	53min 50sec	16min 10sec
Code Minimalism	ResNet	20min 47sec	7min 51sec
Flexibility	Pre-Activation ResNet	20min 57sec	8min 25sec
	DenseNet	49min 23sec	20min 1sec

- PyTorch is quick, but JAX is simply faster
- For small models, the difference can be large
- For large models with few kernel calls, runtimes become similar





Properties	O PyToro
Design	Object-Oriente
Speed	
Community / Doc.	
Debugging	
Code Minimalism	
Flexibility	





- JAX's documentation is good and growing
- Subpackages of JAX (Flax/Optax) sometimes a bit sparser
- PyTorch still has a larger community and more detailed documentation







Properties	Ore PyTorch	
Design	Object-Oriented	Functional Programming
Speed		
Community / Doc.	ÐÐ	Ð
Debugging		
Code Minimalism		
Flexibility		



- For PyTorch, you can print any tensor at any time
- For JAX, JIT-compiled function do not print / debugger difficult
- Once function is jitcompiled, no error can be thrown \Rightarrow NaNs for invalid operations, default behavior for out-ofbounds indexing











Properties	O PyToro
Design	Object-Oriente
Speed	
Community / Doc.	ÐÐ
Debugging	
Code Minimalism	
Flexibility	





- MLP in PyTorch is only a few lines
- PyTorch Lightning reduces code a lot
- More overhead with JAX/Flax, but minor for larger projects





Properties	O PyToro
Design	Object-Oriente
Speed	G
Community / Doc.	ÐÐ
Debugging	
Code Minimalism	
Flexibility	





- JAX: explicit gradient control, removing duplicate tensors, ...
- PyTorch: easy access to tensors, model, etc.
- Both focus on different flexibilities





JAX: Accelerated NumPy with Autograd

• Matrix Operations: JAX shares the same basic API with NumPy

```
import jax
import jax.numpy as jnp
print("Using jax", jax.__version__)
```

Using jax 0.3.25

```
a = jnp.zeros((2, 5), dtype=jnp.float32)
print(a)
```

[[0. 0. 0. 0. 0.]][0. 0. 0. 0. 0.]]





JAX: Accelerated NumPy with Autograd

• Accelerator Support: JAX supports operations on various backends, e.g. GPUs





Tensor by default on accelerator (here GPU) JAX arrays on CPU are NumPy arrays



JAX: Accelerated NumPy with Autograd

• Automatic Differentiation: JAX allows to transform functions, such as taking the gradient of a function via jax.grad

def mse_loss(preds, labels): return ((preds - labels) ** 2).mean() print('Loss', mse_loss(jnp.array([1.0, 2.0]), jnp.array([0.0, 1.5]))

Loss 0.625

mse_grad_fn = jax.grad(mse_loss) print('Loss gradient', mse_grad_fn(jnp.array([1.0, 2.0]),

Loss gradient [1. 0.5]





JAX: Function Transformations

- JAX is based on transformations of pure functions Ο
- For this, JAX lifts functions into an intermediate language called jaxpr

 $exmp_preds = jnp_array([1.0, 2.0])$ exmp_labels = jnp.array([0.0, 1.5])



Pure functions = functions without any side effects, limited to input and outputs



JAX: Function Transformations

- JAX offers several function transformations, most importantly:
 - jax.grad Gradient of a function

jax.jit – Just-In-Time Compilation

- jax.vmap Vectorize function
- jax.pmap Parallelize function on multiple devices





JAX: Just-In-Time Compilation

- Naïve execution of operations can lead to considerable overhead in GPU ↔ CPU communication
- Can we do better?
- Yes, with Just-In-Time Compilation!







JAX: Just-In-Time Compilation

- Just-In-Time (JIT) compilation allows for very efficient code with little effort
- Transforming a function via jax.jit uses XLA (Accelerated Linear Algebra) to compile multiple operations together to improve speed and memory usage

mse_jit = jax.jit(mse_loss)

%timeit mse_loss(x_batch, y_batch).block_until_ready()

434 μs ± 13.3 μs per loop (mean ± std. dev. of 7 runs, 10 00 loops each)

%timeit mse_jit(x_batch, y_batch).block_until_ready() 58.7 μs \pm 1.42 μs per loop (mean \pm std. dev. of 7 runs, 1 0000 loops each)

C Experts





JAX: Just-In-Time Compilation

- Improves execution speed
 (operation fusion, specializing for shapes)
- Improves memory usage (fewer intermediate variables)
- Improves portability (any backend supported in XLA)







JAX: The Sharp Bits

- JAX relies on pure functions, which requires function-centric programming
 - **Immutable Tensors**: in-place operations could have side-effects (JIT may still use in-place ops) \bigcirc
 - **Pseudo Random Numbers**: seed needs to be passed explicit to functions, no global seed variable \bigcirc
- JIT-Compilation is shape specific, need more care with dynamic shapes (e.g., graphs or natural language – use padding)
- Debugging in JIT-compiled functions more difficult since compiled functions can't throw an error \rightarrow potentially undesired side-effects







Summary

- JAX is Accelerated NumPy with Autograd
- Key feature: transformations of pure functions
- Just-In-Time compilation for taking full advantage of accelerators



exmp_preds = jnp.array([1.0, 2.0]) exmp_labels = jnp.array([0.0, 1.5])

jax.make_jaxpr(mse_loss)(exmp_preds, exm

```
{ lambda ; a:f32[2] b:f32[2]. let
    c:f32[2] = sub a b
   d:f32[2] = integer_pow[y=2] c
   e:f32[] = reduce_sum[axes=(0,)] d
   f:f32[] = div e 2.0
 in (f,) }
```





Neural Networks with JAX

- How can we implement an NN in JAX?
- Several libraries available, such as:

Flax – Google Brain, focuses on flexibility and clarity \bigcirc

- Haiku DeepMind, focuses on simplicity and compositionality \bigcirc
- Trax Google Brain, solutions for common training tasks \bigcirc
- Equinox Patrick Kidger and Cristian Garcia, NNs as callable Pytrees \bigcirc







Flax: Neural Networks with JAX

- Main aspects of a Neural Network library in JAX:
 - How can we implement Neural Network layers as functions? \bigcirc
 - How do we handle parameters (weights, biases, etc.)? \bigcirc
 - How do we optimize the model's parameters? \bigcirc
 - How do we put everything together with JIT support? \bigcirc









Flax defines layers as Modules, acting as an immutable dataclass

from flax import linen as nn

class MyModule(nn.Module): # Some dataclass attributes, like hidden dimension # varname : vartype

def setup(self): # Flax uses "lazy" initialization. # In here, define your submodules etc. pass

def ___call__(self, x): # Function for forward pass pass









Flax defines layers as Modules, acting as an immutable dataclass

```
class SimpleClassifier(nn.Module):
    num_hidden : int # Number of hidden neurons
   num_outputs : int # Number of output neurons
   def setup(self):
       # Create the modules we need to build the network
       # nn.Dense is a linear layer
        self.linear1 = nn.Dense(features=self.num_hidden)
        self.linear2 = nn.Dense(features=self.num_outputs)
    def __call__(self, x):
       # Forward pass
       x = self.linear1(x)
       x = nn_tanh(x)
       x = self_linear2(x)
        return x
```

C Experts







We can combine sub-module definition and their call via nn.compact:

```
class SimpleClassifierCompact(nn.Module):
    num_hidden : int # Number of hidden neurons
    num_outputs : int # Number of output neurons
   <u>@nn.compact</u>
   def ___call__(self, x):
        # Forward pass while defining necessary layers
        x = nn.Dense(features=self.num_hidden)(x)
        x = nn_tanh(x)
        x = nn.Dense(features=self.num_outputs)(x)
        return x
```









- Parameters cannot be part of the Module since they must be *mutable*
- Solution: parameters act as an additional input to the module
- Create parameters via module.init function:

example_input = jnp.zeros((16, 2))

model = SimpleClassifier(num_hidden=4, num_outputs=1)









- Parameters are stored as **Immutable Pytrees**, a tree-structured container
 - In Flax, Pytrees are mostly nested dictionaries, with leafs being the parameter values
 - Pytrees allow functions to access collections like parameters, and potentially give one as output
- JAX defines several operations on Pytrees, example: print the shapes of the parameters





- Sub-module names
- Parameters within sub-module



Run a module with parameters via module.apply:

example_input = jax.random.normal(jax.random.PRNGKey(0), (4, 2))

model.apply(params, example_input)

```
DeviceArray([[-0.07297172],
             [ 0.22177655],
             [-0.18521681],
             [-0.165456 ]], dtype=float32)
```







- Several common network layers have been predefined in the Linen API, such as:
 - Linear Layers (nn.Dense, nn.DenseGeneral) \bigcirc
 - Convolutions (nn.Conv, nn.ConvTranspose, etc.) Ο
 - Normalizations (nn.BatchNorm, nn.LayerNorm, etc.) \bigcirc
 - Attention mechanisms (nn.SelfAttention, etc.) \bigcirc
 - Recurrent Neural Networks (nn.LSTMCell, nn.GRUCell, etc.) \bigcirc







Flax: Neural Networks with JAX

- Main aspects of a Neural Network library in JAX:
 - How can we implement Neural Network layers as functions? \bigcirc
 - How do we handle parameters (weights, biases, etc.)? \bigcirc
 - How do we optimize the model's parameters? \bigcirc
 - How do we put everything together with JIT support? \bigcirc









Optax: NN Optimization in JAX

- Optax is a gradient processing and optimization library for JAX based on Pytrees
- Provides building blocks and common optimizers (SGD, Adam, etc.) in a similar function-oriented fashion as Flax

optimizer = optax.adam(learning_rate=1e-3)

opt_state = optimizer.init(params)

updates, opt_state = optimizer.update(grads, opt_sta params)

params = optax.apply_updates(params, update)





- How can we combine the model execution, gradient calculation, and optimizer step, while allowing for Just-In-Time compilation?
- Flax offers a solution with the flax.training sub-library, in particular: TrainState
 - Immutable dataclass with model forward function, parameters, and optimizer (can be extended) \bigcirc

from flax.training.train_state import TrainState model_state = TrainState.create(apply_fn=model.apply, params=params, tx=optimizer)







• A TrainState object can be used as input argument to a function, on which we may apply function transformations (jax.grad, jax.jit, etc.)

Example: binary classification

```
def calculate_loss(state, params, batch):
   data_input, labels = batch
    logits = state.apply_fn(params, data_input)
    logits = logits.squeeze(axis=-1)
    loss = optax.sigmoid_binary_cross_entropy(logits,
    loss = loss.mean()
    return loss
```







• Combine everything into a function that executes a whole training step:

```
@jax.jit
def train_step(state, batch):
    loss, grads = grad_fn(state, state.params, batch)
    state = state.apply_gradients(grads=grads)
    return state, loss
```









• To train the model, we can now just write a training loop that calls the training function several times for different input batches

def train_model(state, data_loader, num_epochs=100): for epoch in range(num_epochs): for batch in data_loader: state, loss = train_step(state, batch) return state

trained_model_state = train_model(model_state, train_data_loader, num_epochs=100)









Flax: Neural Networks with JAX

- Main aspects of a Neural Network library in JAX:
 - How can we implement Neural Network layers as functions? \bigcirc
 - How do we handle parameters (weights, biases, etc.)? \bigcirc
 - How do we optimize the model's parameters? \bigcirc
 - How do we put everything together with JIT support? \bigcirc









Flax: Neural Networks with JAX

- What we haven't discussed yet:
 - Logging can be done with external libraries (e.g., TensorBoard) \bigcirc
 - Flax supports data loading from any other library (e.g., TensorFlow, PyTorch, etc.) \bigcirc
 - Binding parameters to a specific module for easier evaluation \bigcirc
 - Automatically vectorizing and/or parallelizing via jax.vmap and jax.pmap \bigcirc
 - Writing a research code framework for minimal code overhead \bigcirc
 - And much, much more... \bigcirc









Summary

- Flax is a library for NN tools in JAX
- Using immutable dataclasses for more object-oriented programming "feeling"
- Can be combined with several external libraries for optimization, logging, data loading, etc.



def calculate_loss(state, params, batch): data_input, labels = batch

logits = state.apply_fn(params, data_input logits = logits.squeeze(axis=-1)

loss = optax.sigmoid_binary_cross_entropy(

```
loss = loss.mean()
return loss
```

grad_fn = jax.value_and_grad(calculate_loss, argnums=1)

```
@jax.jit
def train_step(state, batch):
    loss, grads = grad_fn(state, state.params,
    state = state.apply_gradients(grads=grads)
    return state, loss
```





When and why to use JAX with Flax?

Benefits

- JAX is extremely fast with Just-In-Time compilation
- Function transformations are powerful tools to easily parallelize and vectorize your code
- Function-oriented programming is helpful in areas like meta-learning, where one needs explicit gradients

Recommendation: if you are doing research and want to get maximum performance out of your code, give JAX a try!



Drawbacks

- The code overhead is usually larger than in other frameworks
- Not as user-friendly / fail-safe as other frameworks
- Handling dynamic shapes can be annoying
- Community still considerably smaller than e.g. TensorFlow or PyTorch





Goals of this talk

- 1) What features are we looking for in an ML/DL framework?
- 2) What is JAX?
- 3) What sets JAX apart from other frameworks?
- 4) How can you train Neural Networks in JAX with Flax?
- 5) Where can I continue my learning journey into JAX with Flax?





Further resources on JAX with Flax

- The <u>JAX</u> and <u>Flax</u> documentations have great introduction tutorials
- If you are interested in seeing JAX with Flax used in practice, and learn new methods in Deep Learning, check out our <u>UvA Deep Learning tutorials</u>!

DEEP LEARNING 1 (JAX+FLAX)

- Tutorial 2 (JAX): Introduction to JAX+Flax
- JAX as NumPy on accelerators
- **H** Function transformations with Jaxpr
- Implementing a Neural Network with Flax
- 🗄 抹 The Fancy Bits 抹
- 🗄 🌂 The Sharp Bits 🌂

< 7

Tutorial 3 (JAX): Activation Functions

Tutorial 4 (JAX): Optimization and Initialization

Tutorial 5 (JAX): Inception, ResNet and DenseNet

Tutorial 6 (JAX): Transformers and Multi-Head Attention

Tutorial 7 (JAX): Graph Neural Networks

Tutorial 2 (JAX): Introduction to JAX+Flax

Status Finished

Filled notebook: Repo View On Github Author: Phillip Lippe

Welcome to our JAX tutorial for the Deep Learning course at the University of Amsterdam! The following notebook is meant to give a short introduction to JAX, including writing and training your own neural networks with Flax. But why should you learn JAX, if there are already so many other deep learning frameworks like PyTorch and TensorFlow? The short answer: because it can be extremely fast. For instance, a small GoogleNet on CIFAR10, which we discuss in detail in Tutorial 5, can be trained in JAX 3x faster than in PyTorch with a similar setup. Note that for larger models, larger batch sizes, or smaller GPUs, a considerably smaller speedup is expected, and the code has not been designed for benchmarking. Nonetheless, JAX enables this speedup by compiling functions and numerical programs for accelerators (GPU/TPU) *just in time*, finding the optimal utilization of the hardware. Frameworks with dynamic computation graphs like PyTorch cannot achieve the same efficiency, since they cannot anticipate the next operations before the user calls them. For example, in an Inception block of GoogleNet, we apply multiple convolutional layers in parallel on the same input. JAX can optimize the execution of this layer by compiling the whole forward pass for the available accelerator and fusing operations where possible, reducing memory access and speeding up execution. In contrast, when calling the first convolutional layer in PyTorch, the framework does not know that multiple convolutions on the same feature map will follow. It sends each

CO Open in Colab





Slides

(personal website)

