# JAX.
# AN INTRODUCTION TO DEEP LEARNING PROGRAMMING PRINCIPLES

SAMUELE PAPA

# WHO AM I?

PhD student @ UvA & NKI

— Real-world applications of Deep Learning

— Scientific modeling

— Medical imaging

Contact: s.papa@uva.nl

# INSPIRATION AND SOURCES

— Phillip Lippe: notebooks, presentations, and more.

— Official JAX documentation.

— Other sources.
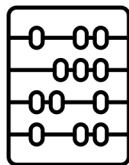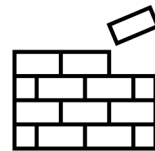
# KEY ASPECTS OF MODERN COMPUTING

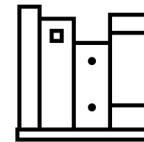# KEY ASPECTS OF MODERN COMPUTING

**Parallelization**

**Compute**

**Memory**

# HOW CAN WE IMPROVE?

## Manual optimization

Use **algorithms** and **data structures**.

e.g. async loading and preprocessing of data on CPU

e.g. hash maps for spatial embeddings

# Leverage compilers

# HOW CAN COMPILERS HELP?

```
[4]   1 import jax
      2 import jax.numpy as jnp
      3
      4 def selu(x, alpha=1.67, lambda_=1.05):
      5   return lambda_ * jnp.where(x > 0, x, alpha * jnp.exp(x) - alpha)
      6
      7 x = jnp.arange(1000000)
      8 %timeit selu(x).block_until_ready()
```

⤵  1.02 ms ± 324 µs per loop (mean ± std. dev. of 7 runs, 1000 loops each)

# HOW CAN COMPILERS HELP?

```
[2]    1 selu_jit = jax.jit(selu)
       2
       3 # Pre-compile the function before timing...
       4 selu_jit(x).block_until_ready()
       5
       6 %timeit selu_jit(x).block_until_ready()
```

264 µs ± 64.4 µs per loop (mean ± std. dev. of 7 runs, 1000 loops each)

# WHY?

Each operation (approximately) calls a new kernel

```
[3]   1 def selu(x, alpha=1.67, lambda_=1.05):
      2     original_x = x
      3     x = jnp.exp(x)
      4     x = alpha * x
      5     x = x - alpha
      6     x = jnp.where(original_x > 0, original_x, x)
      7     x = lambda_ * x
      8     return x
      9
     10 x = jnp.arange(1000000)
     11 %timeit selu(x).block_until_ready()
```

    922 μs ± 114 μs per loop (mean ± std. dev. of 7 runs, 1000 loops each)

# JUST-IN-TIME COMPILATION

**JIT. Compile code during execution.**

**Simply use** `jax.jit()`

Compiles the function by converting to intermediate `jaxprs` language.

Tracks usage and optimizes also memory.

# FUNCTIONAL PROGRAMMING

# PURE FUNCTIONS

Example of **classic object-oriented** design pattern.
Often encountered when using classes.

```
[4]    1 counter = 0
       2 def increase_counter_by(x):
       3     return counter + x
       4
       5 print(increase_counter_by(12))
```

12

```
[5]    1 counter = 10
       2 print(increase_counter_by(12))
```

22

# PURE FUNCTIONS

Now we compile it.

```
[6]   1 jit_increase_counter_by = jax.jit(increase_counter_by)
      2 print(jit_increase_counter_by(10))

      20


[7]   1 counter = 0
      2 print(jit_increase_counter_by(12))

      22
```

# PURE FUNCTIONS

A pure function is a function that, given the **same input**, will always return the **same output** and does not have any observable **side effect**.

*A side effect is e.g. something that is performed in-place, affects something outside the scope of the function.*

# WHY PURE FUNCTIONS?

1.     Makes your code more maintainable.

2.     Makes compilation possible and simple.

3.     Makes parallelization easier.

4.     You can replace the whole function with its outputs when necessary.

5.     Functional composition makes math-to-code easier.

# HOW TO WRITE JAX

# JAX IS NUMPY

array concept, just as in numpy.

```
[9]    1 import jax.numpy as jnp
       2 a = jnp.zeros((2, 5), dtype=jnp.float32)
       3 print(a)
```

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

All numpy functions are available.

API matches.

```
[10]   1 b = jnp.arange(6)
       2 print(b)
```

```
   [0 1 2 3 4 5]
```

# JAX IS NUMPY

`Array` objects are always placed directly on the available accelerators

When we retrieve from device, it becomes a `numpy` array.

```
[13]    1 b.devices()

    {cuda(id=0)}
```

```
[14]    1 b_cpu = jax.device_get(b)
        2 print(b_cpu.__class__)

    <class 'numpy.ndarray'>
```

# PARALLELIZATION

# JAX HAS AUTOMATIC VECTORIZATION

Simple parallelization of operations
using `jax.vmap()`

An example of simple element-wise operation.

```
[13]   1 def single_linear(x, w, b):
       2     return (x[:,None] * w).sum(axis=0) + b
       3
       4 # Example inputs
       5 x_in = jnp.ones((4,))
       6 w_in = jnp.ones((4, 3))
       7 b_in = jnp.ones((3,))
       8
       9 single_linear(x_in, w_in, b_in)
```

```
Array([5., 5., 5.], dtype=float32)
```

# JAX HAS AUTOMATIC VECTORIZATION

**Input shapes**

[5,4] – batched

[4,3] – shared

[4]    – shared

```
[14]   1 batched_linear = jax.vmap(
       2       single_linear,
       3       in_axes=(0,None,None),
       4       out_axes=0
       5 )
       6
       7 # Example batched inputs
       8 batched_x_in = jnp.ones((5, 4,))
       9
      10 batched_linear(batched_x_in, w_in, b_in)
```

**Output shapes**

[4,3] – batched

```
Array([[5., 5., 5.],
       [5., 5., 5.],
       [5., 5., 5.],
       [5., 5., 5.],
       [5., 5., 5.]], dtype=float32)
```

# PARALLELIZATION

# FUNCTIONAL COMPUTATION OF GRADIENTS

The `jax.grad()` function returns the function that evaluates the derivative at any given input

```python
[15]    1 def rms_error(x, y):
        2     return jnp.sqrt(jnp.mean((x-y)**2))
        3
        4 x_in = jnp.array([1.2,3.2,4], dtype=jnp.float32)
        5 y_target = jnp.array([0,5.,10.], dtype=jnp.float32)
        6
        7 grad = jax.grad(rms_error)
        8 grad(x_in, y_target)
```

    Array([ 0.1086251 , -0.16293764, -0.54312545], dtype=float32)

# GRADIENT DESCENT

Very intuitive implementation
from math to code.

```python
[22]  1 lambda_ = 0.05
      2 x_new = x_in
      3 for i in range(200):
      4     x_new = x_new - lambda_ * grad(x_new, y_target)
      5     if i % 10 == 0:
      6         print(rms_error(x_new, y_target))
```

```
3.6657236
3.4990568
3.3323894
3.1657221
2.9990551
2.8323882
2.665721
2.4990537
2.3323865
2.1657195
1.9990524
1.8323854
1.6657186
1.4990516
1.3323847
1.1657186
0.9990542
0.8323898
0.6657253
0.49906078
```

```python
[23]  1 x_new
```

```
Array([0.11374928, 4.829371  , 9.431246  ], dtype=float32)
```

# GRADIENT DESCENT
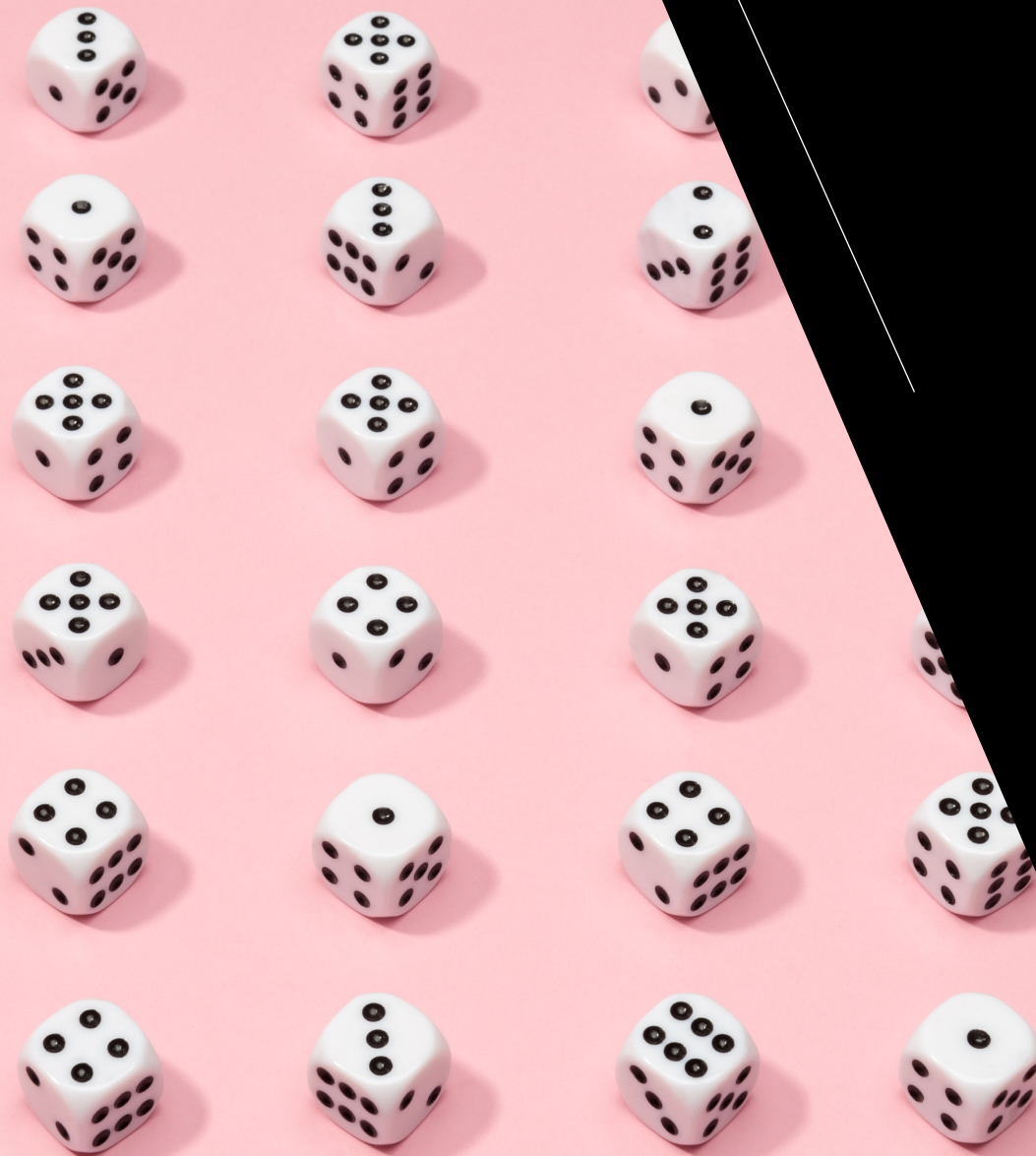
Improving the performance a bit
with some heuristic annealing

```
[18]   1 lambda_ = 0.1
       2 annealing = 0.75
       3 steps = 50
       4 x_new = x_in
       5 for i in range(1000):
       6     x_new = x_new - lambda_ * grad(x_new, y_target)
       7     if (i + 1) % steps == 0:
       8         lambda_ *= annealing
       9         print(rms_error(x_new, y_target))
```

```
2.0157194
0.7657221
0.015718736
0.012406095
0.008687421
0.0071328944
0.0047321706
0.0041667605
0.002507655
0.0024983105
0.0012557198
0.0012558495
0.0008558435
0.0007279681
0.00046012658
0.00043075677
0.00023763097
0.00023762452
0.00013843694
0.00013850731
```

```
[19]   1 x_new
```

```
Array([4.142728e-05, 4.999917e+00, 9.999779e+00], dtype=float32)
```
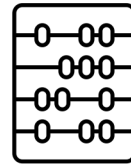
# PSEUDO-RANDOM NUMBER GENERATION
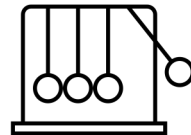
# THE GOALS OF PSEUDO-RANDOM NUMBER GENERATOR

**Reproducible**

**Parallelizable**

**Vectorizable**

# HOW IT IS USUALLY DONE

**Set a global seed.**

How does it behave on multiple devices?

What happens with intermediate steps of random sampling?

```
[21]    1 import numpy as np
        2 import torch
        3 np.random.seed(0)
        4 torch.manual_seed(0)
```

⤳ <torch._C.Generator at 0x7f32c06c33d0>

# WHEN DOES THE GLOBAL SEED FAIL?

**Order of operations is not guaranteed.**

Especially in parallel computations.

```python
[22]  1 import numpy as np
      2
      3 np.random.seed(0)
      4
      5 def bar(): return np.random.uniform()
      6 def baz(): return np.random.uniform()
      7
      8 def foo(): return bar() + 2 * baz()
      9
     10 print(foo())
```

1.9791922366721637

# WHEN DOES THE GLOBAL SEED FAIL?

```python
[26]    1 import numpy as np
        2
        3 np.random.seed(0)
        4
        5 def bar(): return np.random.normal()
        6 def baz(): return np.random.uniform()
        7
        8 def foo(): return bar() + baz()
        9
       10 print(foo())
```
⮏ 2.366815722039308

```python
[27]    1 import numpy as np
        2
        3 np.random.seed(0)
        4
        5 def bar(): return np.random.uniform()
        6 def baz(): return np.random.normal()
        7
        8 def foo(): return bar() + baz()
        9
       10 print(foo())
```
⮏ 1.290405244736486

**Same operation, different results.**

# USE PRNG KEYS

Key: used by pseudo-random number generator to actually create randomness.

Given a key, the output of the random operation is always the same.

```
[28]  1 from jax import random
      2
      3 key = random.key(42)
      4 print(key)
```

Array((), dtype=key<fry>) overlaying:
[ 0 42]

```
[29]  1 print(random.uniform(key))
      2 print(random.uniform(key))
```

0.42672753
0.42672753

Same is possible in `numpy` and `torch` using generators.

# DID WE SOLVE THE PROBLEM?

```
[31]   1 key = random.key(42)
       2
       3 def bar(key): return random.uniform(key)
       4 def baz(key): return random.normal(key)
       5
       6 def foo(key): return bar(key) + baz(key)
       7
       8 print(foo(key))
```

⇥ 0.24201576

```
[32]   1 key = random.key(42)
       2
       3 def bar(key): return random.normal(key)
       4 def baz(key): return random.uniform(key)
       5
       6 def foo(key): return bar(key) + baz(key)
       7
       8 print(foo(key))
```
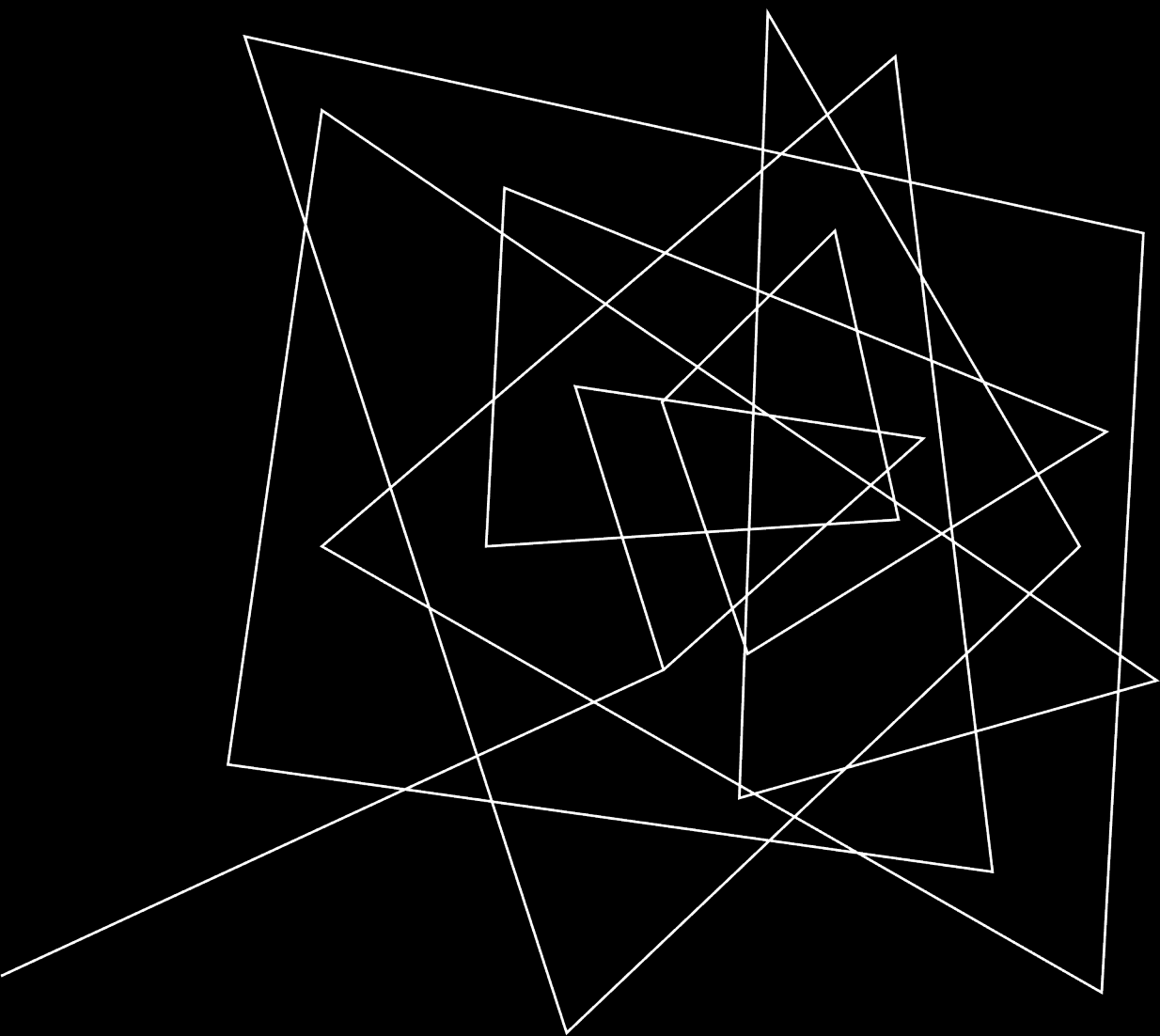
⇥ 0.24201576

# SUMMARY

1. **Compilation** = free code optimization.

2. **Functional** programming is powerful.

3. **Vectorization** to explicitly batch operations.

4. JAX at the core is numpy with **autograd.**

5. Reliable **pseudo-RNG** with keys.

# THANK YOU!

Samuele Papa.
Open to chat and collaborate!
Looking for internship opportunities.

samuele.papa@gmail.com
samuelepapa.github.io