# JAX - GPU accelerated linear algebra & Wigner's Batman

Hrvoje Abraham

# JAX by Google

Originally, it stood for "Just After eXecution"



Transformable numerical computing at scale

JAX is a Python library for accelerator-oriented array computation and program transformation, designed for high-performance numerical computing and large-scale machine learning.

# JAX by Google

Type '/' to search projects

Help     Docs     Sponsors     Log in     Register

jax 0.8.1

✓  Latest version

pip install jax

Released: Nov 18, 2025

Differentiate, compile, and transform Numpy code.

https://pypi.org/project/jax/

https://github.com/jax-ml/jax

```python
import jax
import jax.numpy as jnp

def predict(params, inputs):
  for W, b in params:
    outputs = jnp.dot(inputs, W) + b
    inputs = jnp.tanh(outputs)  # inputs to the next layer
  return outputs                # no activation on last layer

def loss(params, inputs, targets):
  preds = predict(params, inputs)
  return jnp.sum((preds - targets)**2)

grad_loss = jax.jit(jax.grad(loss))  # compiled gradient evaluation function
perex_grads = jax.jit(jax.vmap(grad_loss, in_axes=(None, 0, 0)))  # fast per-example grads
```

# JAX by Google

| Hardware | Backend | Support Status |
|---|---|---|
| **NVIDIA GPUs** | CUDA | First-class, official support. |
| **Google TPUs** | TPU | First-class, official support. |
| **AMD GPUs** | ROCm | Official support (Linux only). |
| **Apple Silicon** | Metal | Experimental (via `jax-metal` plugin). |
| **CPU** | Host | Fallback for all systems. |

No direct Vulkan support.

# JAX by Google

| | Linux, x86_64 | Linux, aarch64 | Mac, aarch64 | Windows, x86_64 | Windows WSL2, x86_64 |
|---|---|---|---|---|---|
| CPU | yes | yes | yes | yes | yes |
| NVIDIA GPU | yes | yes | n/a | no | experimental |
| Google Cloud TPU | yes | n/a | n/a | n/a | n/a |
| AMD GPU | yes | no | n/a | no | experimental |
| Apple GPU | n/a | no | experimental | n/a | n/a |
| Intel GPU | experimental | n/a | n/a | no | no |

# Example 1 – Matrix multiplication

**Matrix 1**

| 1 | 1 |
|---|---|
| 2 | 2 |
| 3 | 3 |

(3x2)

**X**

**Matrix 2**

| 1 | 1 | 1 |
|---|---|---|
| 2 | 2 | 2 |

(2x3)

**=**

**Product matrix**

| 3 | 3 | 3 |
|---|---|---|
| 6 | 6 | 6 |
| 9 | 9 | 9 |

(3x3)

$$c_{ij} = \text{row } i \times \text{col } j = \sum_{k=1}^{n} a_{ik} b_{kj}$$

$$\text{row } i \begin{bmatrix} \rule{2cm}{0.4pt} \\ \\ \end{bmatrix} \times \overset{\text{col } j}{\begin{bmatrix} \\ \\ \end{bmatrix}} = \begin{bmatrix} \square \\ \\ \end{bmatrix}$$

$$\underset{m \times n}{A} \qquad \times \qquad \underset{n \times p}{B} \qquad = \qquad \underset{m \times p}{C}$$

Example 1 – Matrix multiplication



# Nvidia RTX 4080 16GB VRAM    Na=30,000 x Nb=30,000 => Nc=30,000

# Example 1 – Matrix multiplication

```python
MATRIX_SIZE = 32000
DTYPE = jnp.float32
```

```python
A = jax.random.normal(k1, (N, N), dtype=dtype)
B = jax.random.normal(k2, (N, N), dtype=dtype)
```

```python
total_sum = jnp.sum(jnp.matmul(A, B))
```

```python
flops = 2 * (N ** 3)
tflops = flops / duration / 1e12

print(f"Total Sum: {total_sum}")
print(f"Execution time: {duration:.4f} s")
print(f"Performance:    {tflops:.2f} TFLOPS")
```



```
Cmd : 0    △ WSL : 15    +         default : WSL:Ubuntu-24.04   C Z N

c/repos/pinn → python3 jax_test_4.py
```

# Example 1 – Matrix multiplication

RTX 4080 16GB VRAM    Na=20,000 x Nb=20,000 => Nc=20,000

```
⚡  ▣Cmd : 0   △ WSL : 15   +        default : WSL:Ubuntu-24.04  ℂ ℤ ℕ  08:01:30  🕐 Dec 09 15:41:30

c/repos/pinn → python3 jax_test_4.py                    at 23:42:26  took 1s778ms
```

0.33s @ 48 TFLOPs

16,000 billion operations

# Example 1 – Matrix multiplication

| Device | Approx. Peak Performance (FP32 TFLOPS) | Time for 50 TFLOPs |
|---|---|---|
| NVIDIA RTX 4080 GPU | ~48.7 TFLOPS  TechPowerUp +1 | ~1.03 seconds |
| Intel Core i9-13800H CPU | ~1.5–2.0 TFLOPS (estimated from benchmarks  NanoReview +1 ) | ~25–33 seconds |
| HP Notebook (5 years old, midrange) | ~0.2–0.3 TFLOPS (typical integrated GPU or older CPU  Microsoft Learn +1 ) | ~3–4 minutes |

**200x speedup**

# Example 1 – Matrix multiplication – Warmup phase

# Example 1 – Matrix multiplication – Warmup phase

# Example 1 – Matrix multiplication – Warmup phase

# Example 1 – Matrix multiplication – Warmup phase



https://openxla.org/xla

XLA 🔖 ▾

SEND FEEDBACK

XLA (Accelerated Linear Algebra) is an open-source compiler for machine learning. The XLA compiler takes models from popular frameworks such as PyTorch, TensorFlow, and JAX, and optimizes the models for high-performance execution across different hardware platforms including GPUs, CPUs, and ML accelerators.

As a part of the OpenXLA project, XLA is built collaboratively by industry-leading ML hardware and software companies, including Alibaba, Amazon Web Services, AMD, Apple, Arm, Google, Intel, Meta, and NVIDIA.

# Example 1 – Warmup Phase

**Computational structure (code)**



**+**

**Hardware (10,000 GPUs!)**



**=**

**Computational Plan**



Declarative! No need for low-level kernels.

High performance reruns for different inputs.

# Example 2 – Matrix eigenvalues

# Example 2 – Wigner's theorem

(Wigner 1955, 1958). This law was first observed by Wigner (1955) for certain special classes of random matrices arising in quantum mechanical investigations.



The distribution of eigenvalues of a symmetric random matrix with entries chosen from a standard normal distribution is illustrated above for a random $5000 \times 5000$ matrix.

https://mathworld.wolfram.com/WignersSemicircleLaw.html

# Example 2 – Wigner's theorem

matrix element mean = 0, variance = 1/n, std_dev = 1/sqrt(n)

```python
def generate_symmetric_matrix(key, n):
    """

    Generates a symmetric matrix where elements have mean 0 and variance 1/n.
    """

    key, subkey = random.split(key)
    # Generate random normal values with std dev = 1/sqrt(n) -> variance = 1/n
    scale = 1.0 / jnp.sqrt(n)
    A = random.normal(subkey, (n, n)) * scale

    # Symmetrize: Use upper triangle for both upper and lower parts
    # M_ij = A_ij for i <= j
    # M_ji = A_ij for i < j
    M = jnp.triu(A) + jnp.triu(A, 1).T
    return M
```

# Example 2 – Wigner's theorem

```python
print("Calculating eigenvalues...")
start_time = time.time()

# eigh is optimized for Hermitian/symmetric matrices
eigens_jax = jnp.linalg.eigh(M)[0]

# Block to ensure calculation is done
eigens_jax.block_until_ready()
end_time = time.time()
```

# Example 2 – Wigner's theorem, N=15000 in 30s

# Example 2 – Matrix eigenvalues, Wigner's theorem



Eigenvalue Density Distribution (N=15000)

# Example 3 – Wigner's Batman

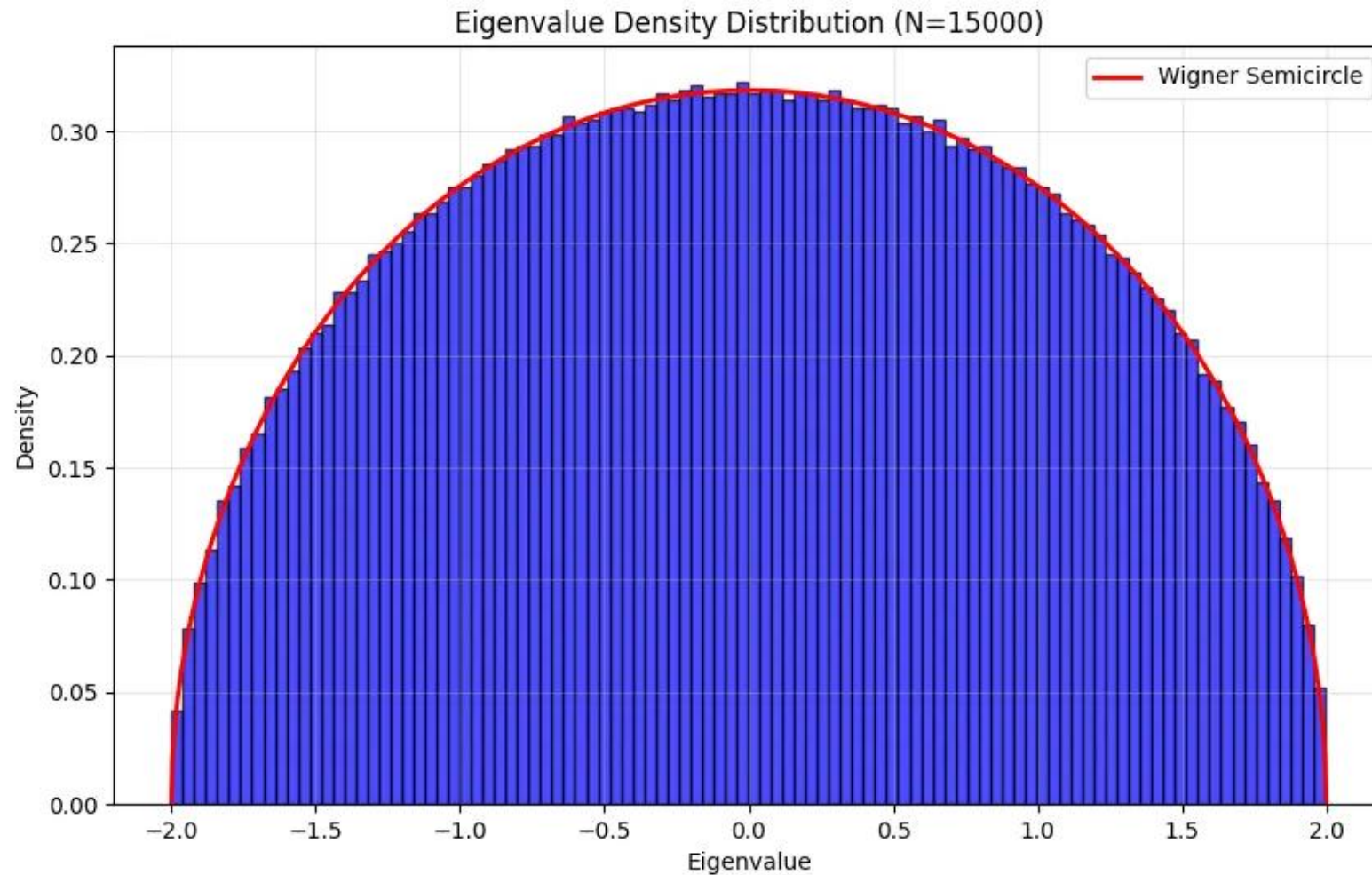## Empirical deviations of semicircle law in mixed-matrix ensembles

Mehmet Süzen*
(Dated: December 3, 2021)

An algorithm is introduced for sampling a set of matrices from mixed orders random matrix ensembles, i.e., Mixed Matrix Ensemble Sampling (MMES). The concept of *the degree of mixture* of the matrix ensemble provides a balanced sampling of the mixed matrix ensemble. As an application of MMES, we have shown how the semicircle law deviates from the conventional behaviour in mixed Gaussian Orthogonal Ensemble (mGOE) as a novel finding.

https://hal.science/hal-03464130v1/document

# Example 3 – Wigner's Batman

```python
# 1. Sample Matrix Orders
matrix_orders = random.binomial(key_orders, N_MAX, MU, shape=(N_MATRICES,))
matrix_orders = jnp.clip(matrix_orders, a_min=2, a_max=N_MAX)
orders_np = np.array(matrix_orders)
```

```python
print("Generating mixed ensemble...")
start_time = time.time()

for i, order in enumerate(orders_np):
    eigs = get_periodic_eigenvalues_sorted(keys[i], int(order), N_MAX)
    eigs.block_until_ready()
    all_eigenvalues.append(eigs)

raw_eigens = jnp.concatenate(all_eigenvalues).tolist()

# 2. Symmetrization
symmetric_eigens = raw_eigens + [-x for x in raw_eigens]

print(f"Done in {time.time() - start_time:.4f}s")
```
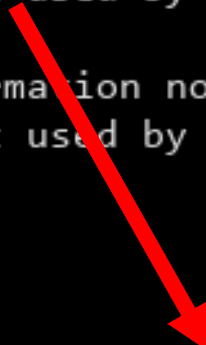
# Example 3 – Wigner's Batman, 100 matrices ensemble in 29s

# Example 3 – Wigner's Batman

Combined distribution of eigenvalues for binomially distributed size matrix ensemble:
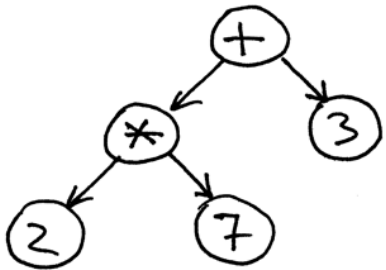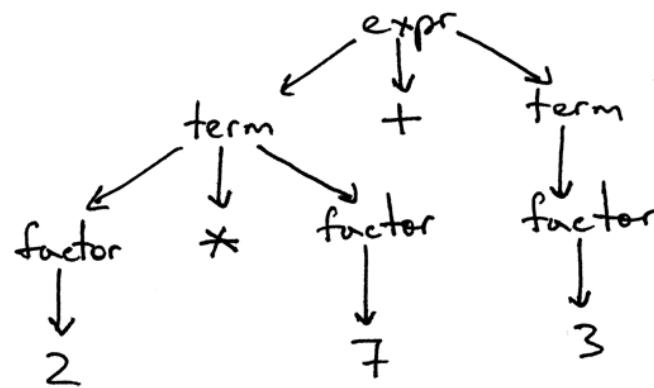
# Example 4 – Compound computation, GPU-focused

**The Task:** $x_{t+1} = \tanh(W \cdot x_t + b) + \alpha \cdot x_t$ **Steps:** $100{,}000$ **Matrix Size:** $2048 \times 2048$

$2 * 7 + 3$

AST

Parse tree

Code is converted into a computational graph which is than transformed and optimized for target platform and final output!!!
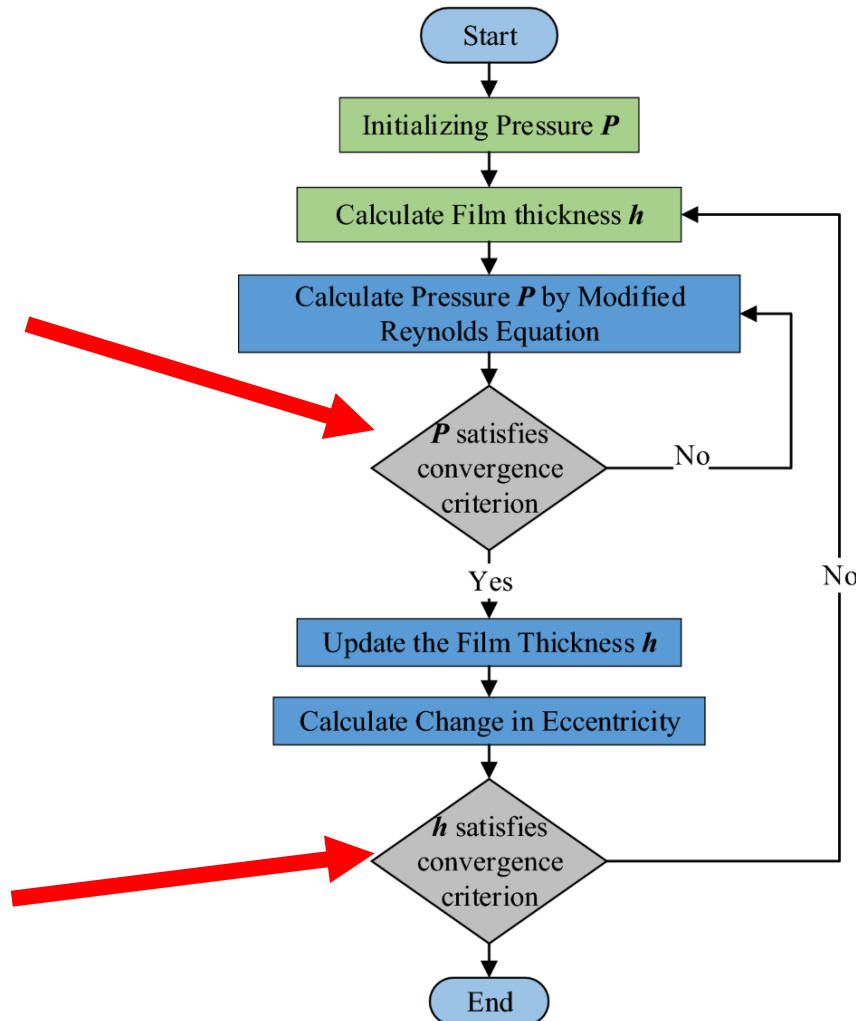
Loops, derivatives, summation…
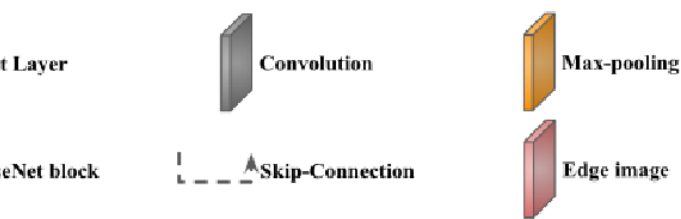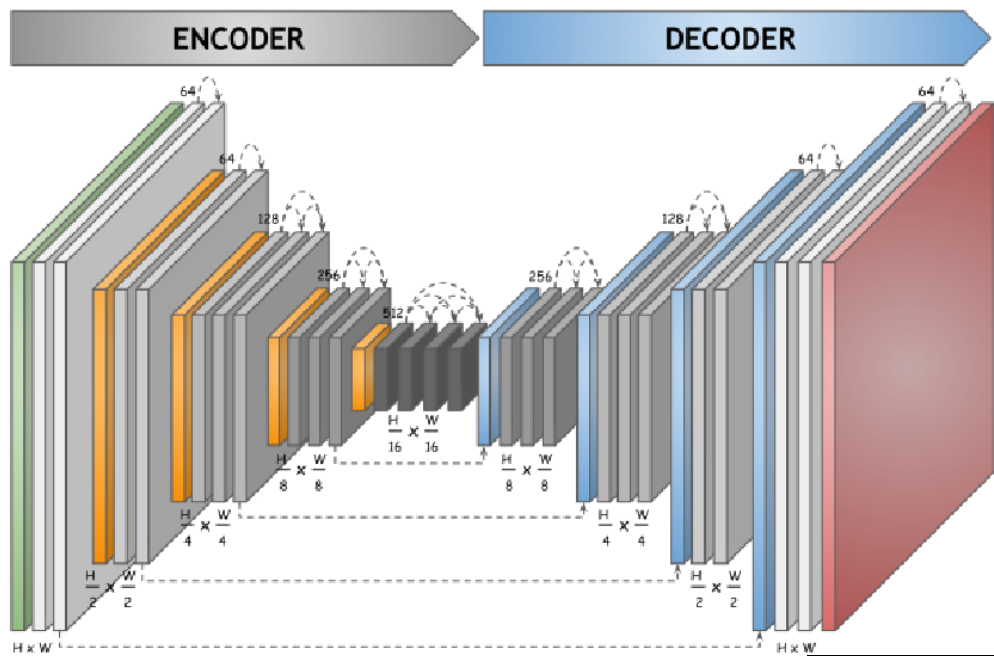
# Example 5 – Dynamic loops, conditional stop, GPU-focused



- Use: solvers, training…

# JAX vs PyTorch

- Problem for high-level pytorch
  - CPU-bound, or:
  - Need to go into low-level *Triton*

- Possible via high-level JAX

# Example 6 – Derivatives, Hessian, AI, Tensor Networks, QC

# Thank you

**AVL** ✦

www.avl.com