



UNIVERSITY OF  
WATERLOO

NengoFPGA

ABR  
APPLIED BRAIN RESEARCH

# NENGOFPGA: AN FPGA BACKEND FOR THE NENGO NEURAL SIMULATOR

Ben Morcos  
August 9, 2019

## — Disclaimer

---

This presentation and work described herein is heavily based on my FPT2018 paper\*.

---

\* Implementing NEF Neural Networks on Embedded FPGAs

## — Claim —

---

We develop a Python accessible PYNQ implementation that outperforms the Jetson TX1 GPU computing NEF style neural networks by  $10\text{--}500\times$  using  $2\text{--}10\times$  less power.

## — Claim —

---

We develop a Python accessible PYNQ implementation that outperforms the Jetson TX1 GPU computing NEF style neural networks by  $10\text{--}500\times$  using  $2\text{--}10\times$  less power.

Our design:

- ▶ Is built for embedded applications
- ▶ Is built with High-Level Synthesis
- ▶ Automatically tunes fixed-point precision

## Motivation

---

- ▶ Machine learning is growing

## Motivation

---

- ▶ Machine learning is growing
  - ↳ Typically done in Python



## Motivation

---

- ▶ Machine learning is growing
  - ↳ Typically done in Python
  - ↳ Compute heavy



## Motivation

---

- ▶ Machine learning is growing
  - ↳ Typically done in Python
  - ↳ Compute heavy
- ▶ Moore's law is ending



## Motivation

---

- ▶ Machine learning is growing
  - ↳ Typically done in Python
  - ↳ Compute heavy
- ▶ Moore's law is ending
  - ↳ Want efficient/custom hardware



## — Context —

---

We target the Neural Engineering Framework (NEF) via the Nengo Python package

We target the Neural Engineering Framework (NEF) via the Nengo Python package

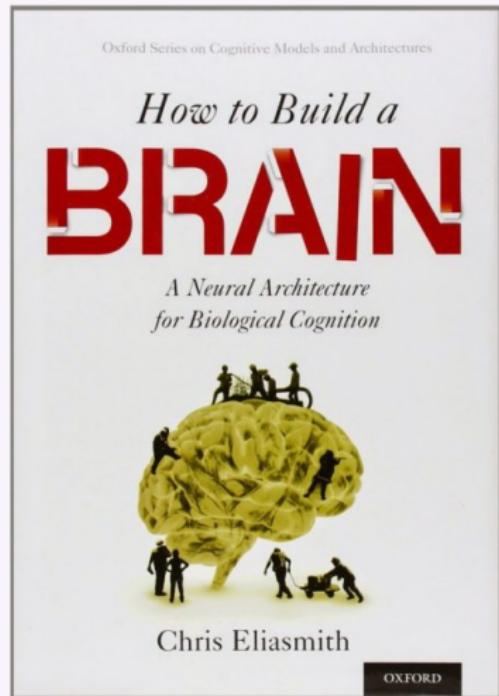
- ▶ Note NEF has different structure than DNN

## — Context —

### Neural Engineering Framework

Nengo goes beyond DNN

- ▶ Dynamic models

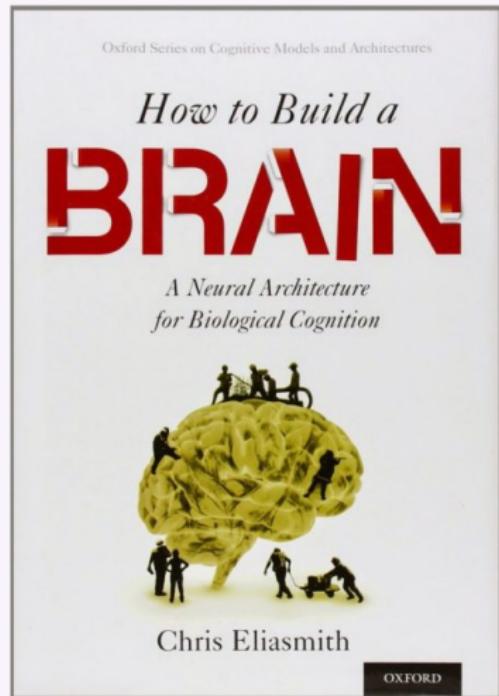


## — Context —

### Neural Engineering Framework

Nengo goes beyond DNN

- ▶ Dynamic models
- ▶ Real-time adaptive control

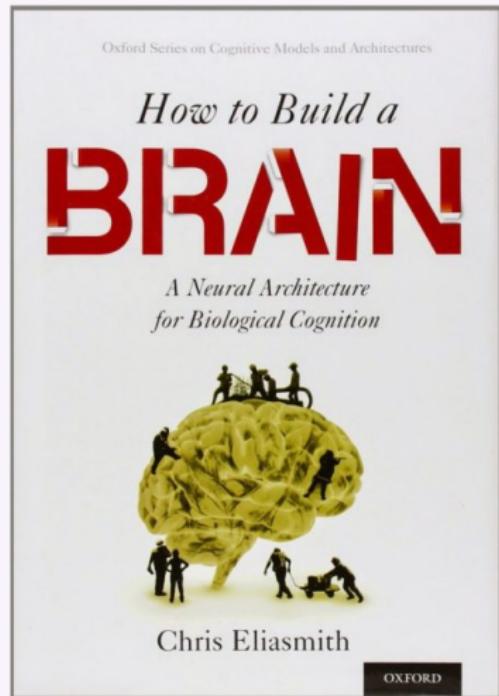


## — Context —

### Neural Engineering Framework

Nengo goes beyond DNN

- ▶ Dynamic models
- ▶ Real-time adaptive control
- ▶ Biologically plausible cognitive architectures

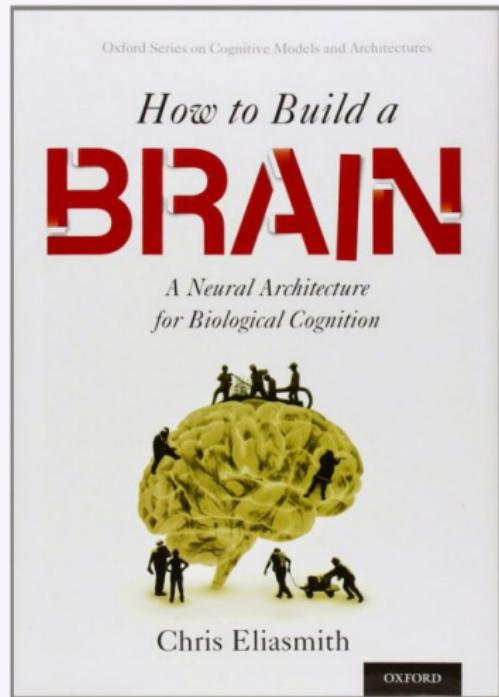


## Context

### Neural Engineering Framework

Nengo goes beyond DNN

- ▶ Dynamic models
- ▶ Real-time adaptive control
- ▶ Biologically plausible cognitive architectures
- ▶ Hierarchical reinforcement learning

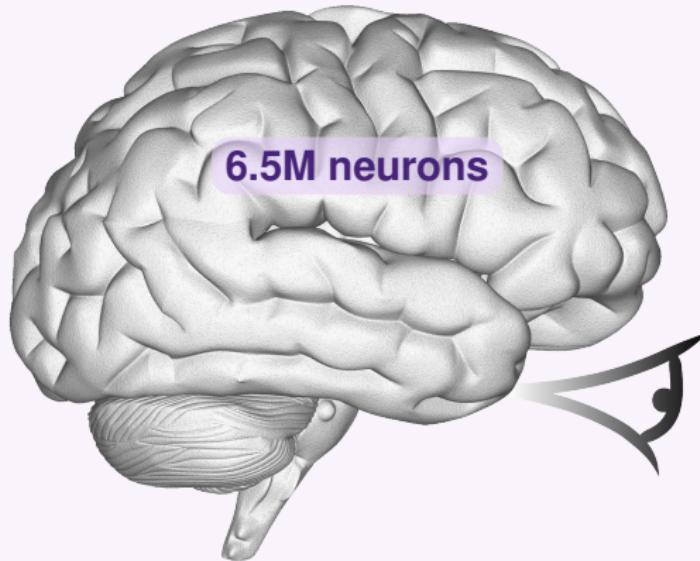


## Context

---

### Neural Engineering Framework

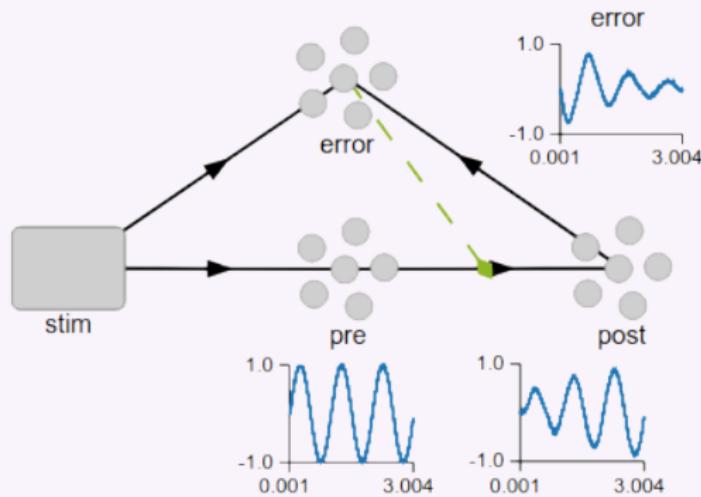
The world's largest functional brain model, Spaun, is built with the NEF and Nengo



## Context

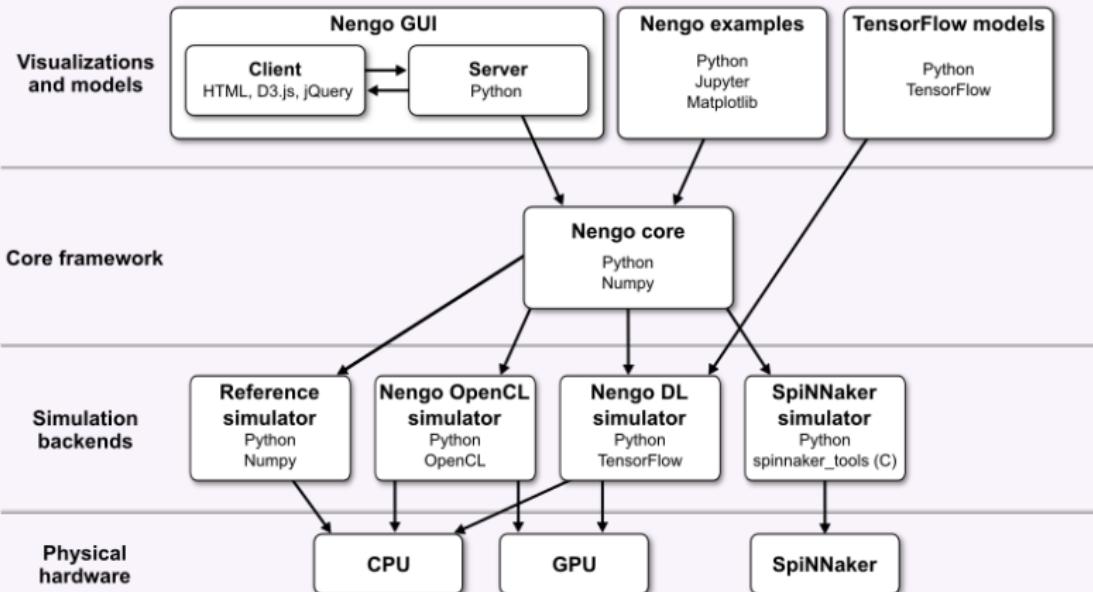
### Nengo Visualization

Online learning displayed with Nengo GUI



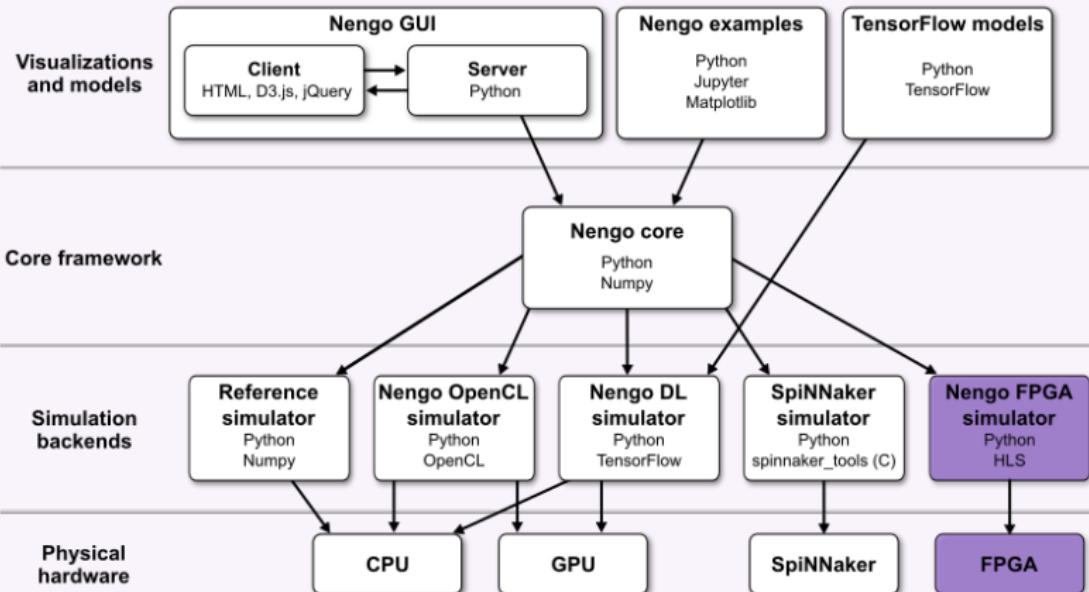
# Context

## Nengo Ecosystem



# Context

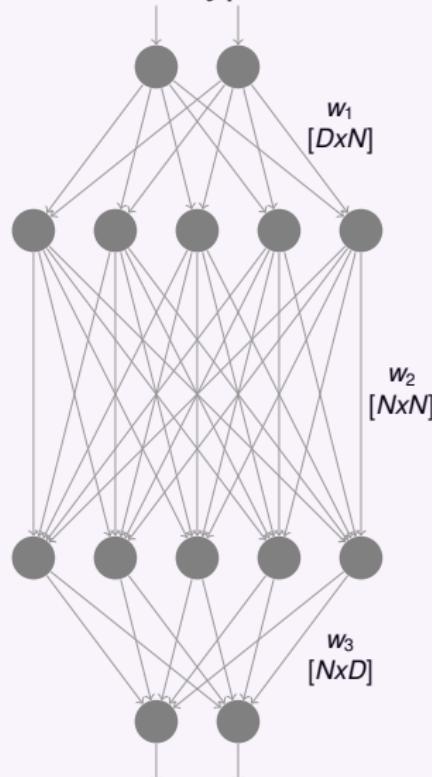
## Nengo Ecosystem



## Context

Neural Engineering Framework

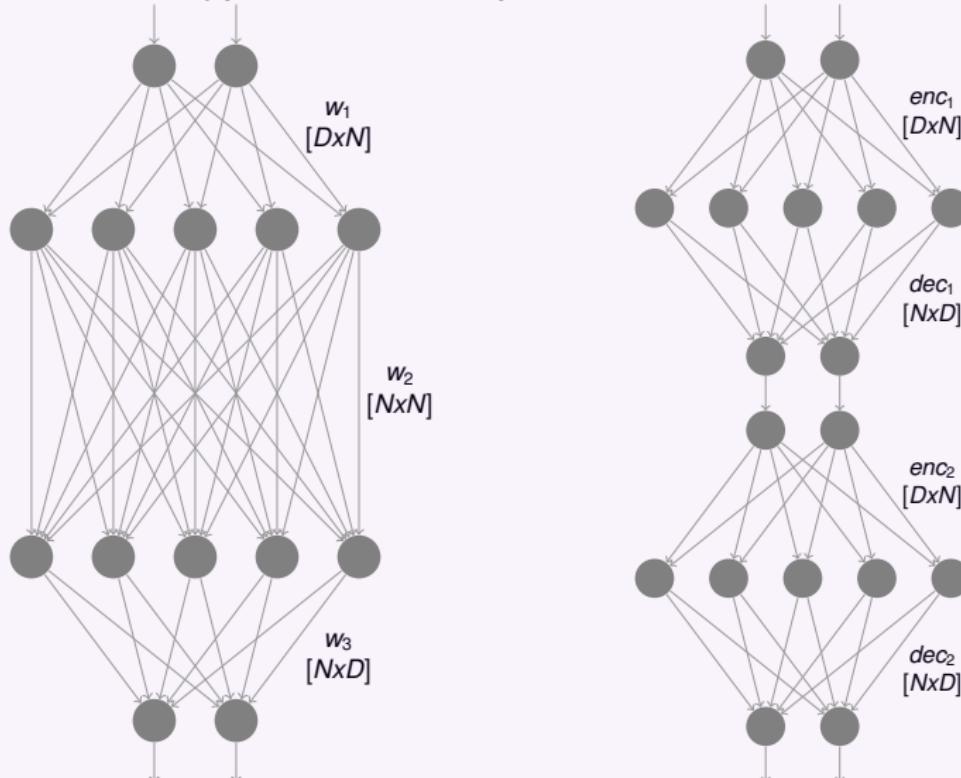
First let's look at typical DNN



## Context

Neural Engineering Framework

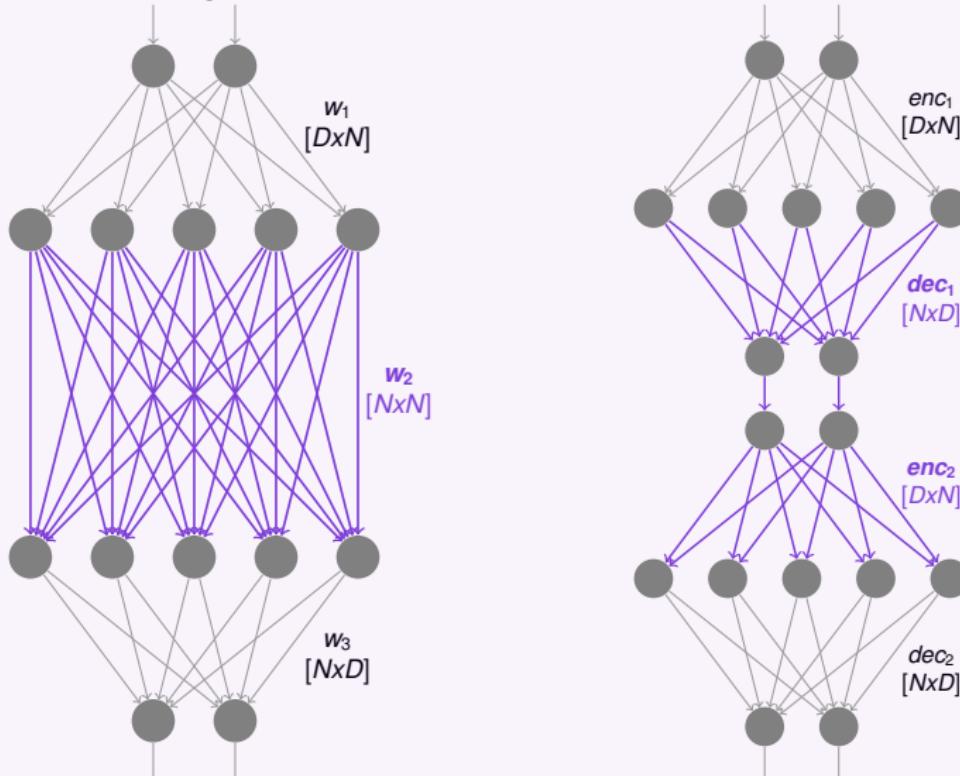
First let's look at typical DNN compared to an NEF network



## Context

Neural Engineering Framework

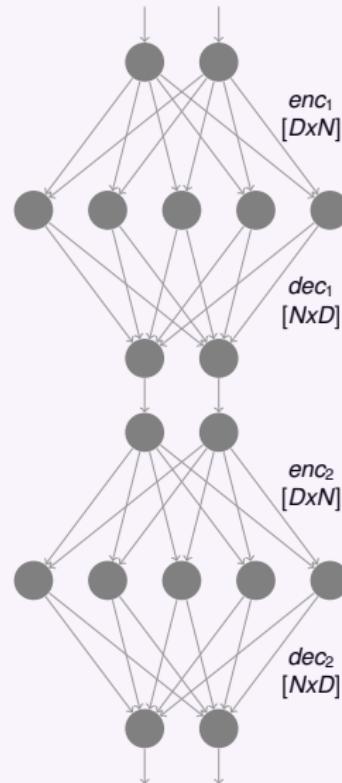
We factor the weight matrix into encode and decode matrices



# Context

## Neural Engineering Framework

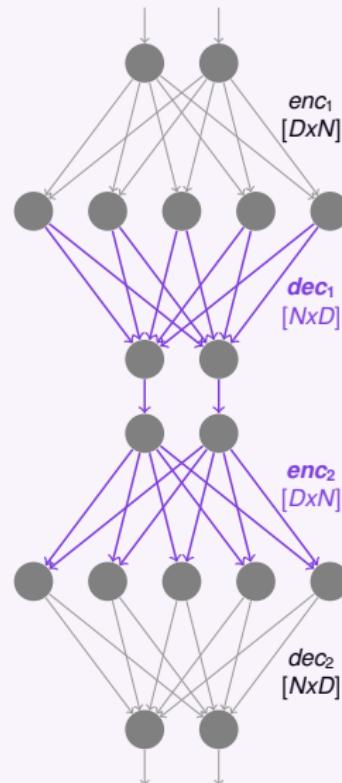
- ▶ We have  $N \gg D$



# Context

## Neural Engineering Framework

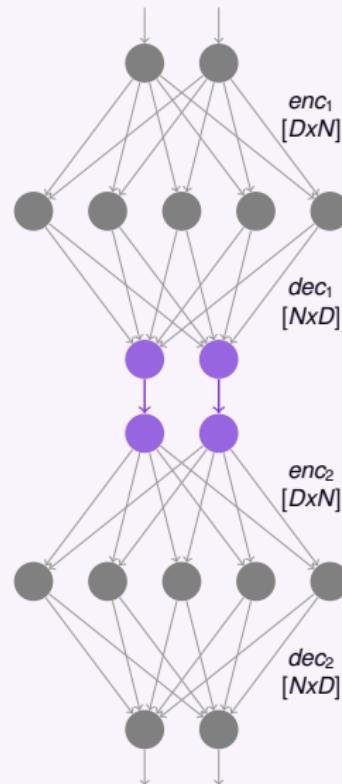
- ▶ We have  $N \gg D$ 
  - ↳ NEF saves memory



# Context

## Neural Engineering Framework

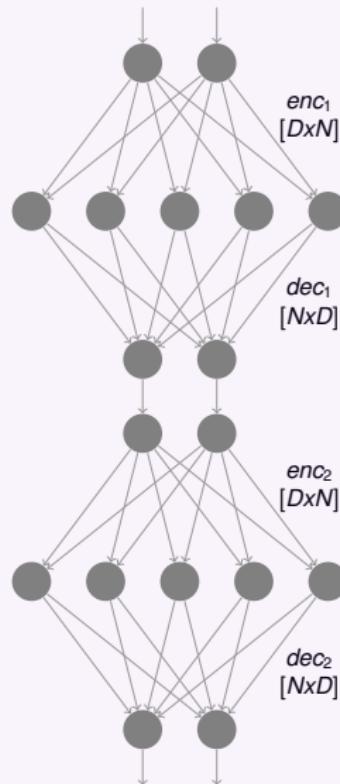
- ▶ We have  $N \gg D$ 
  - ↳ NEF saves memory
  - ↳ NEF saves bandwidth



# Context

## Neural Engineering Framework

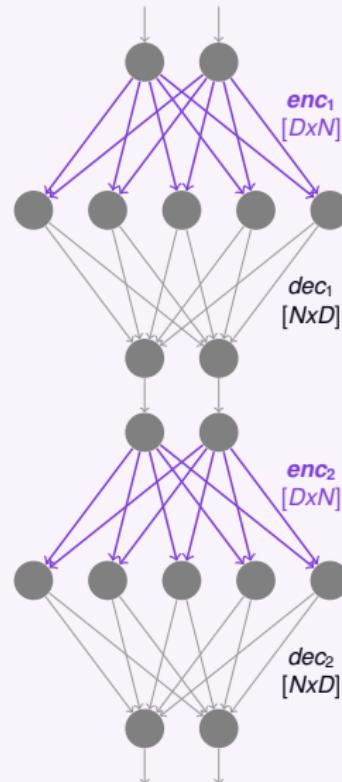
- ▶ We have  $N \gg D$ 
  - ↳ NEF saves memory
  - ↳ NEF saves bandwidth
- ▶ NEF adds temporal dimension



# Context

Neural Engineering Framework

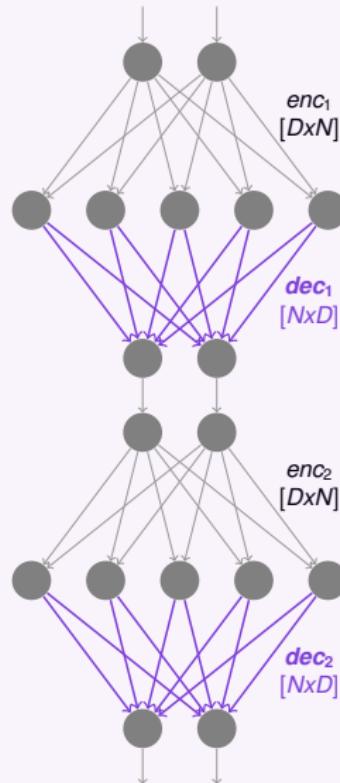
- ▶ Encoders are generated **randomly**



# Context

## Neural Engineering Framework

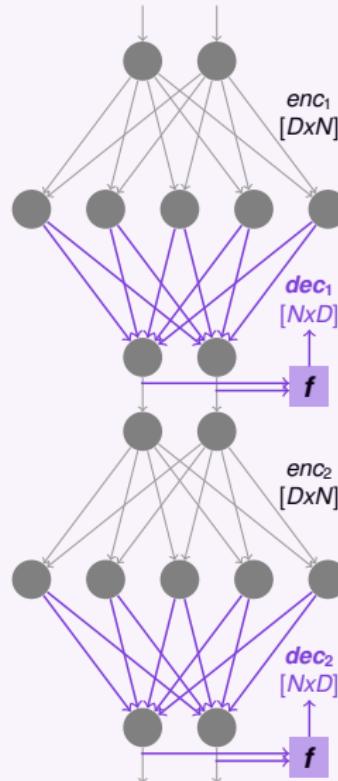
- ▶ Encoders are generated **randomly**
- ▶ Decoders can be:
  - ↪ Pre-solved in offline training and/or



## — Context —

### Neural Engineering Framework

- ▶ Encoders are generated **randomly**
- ▶ Decoders can be:
  - ↪ Pre-solved in offline training and/or
  - ↪ Updated with **online learning**

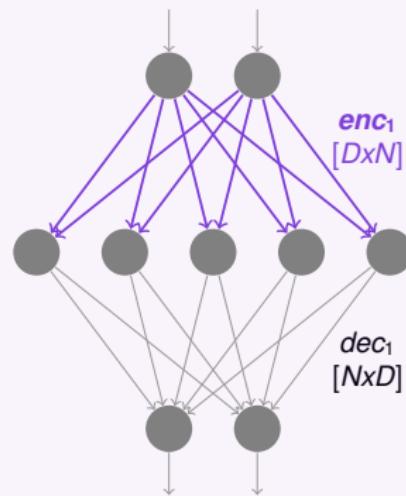
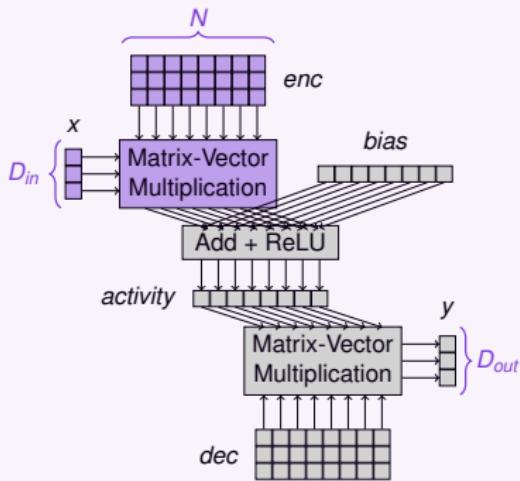


# Context

## Neural Engineering Framework

This boils down to two matrix multiplies

- ▶ One for encode

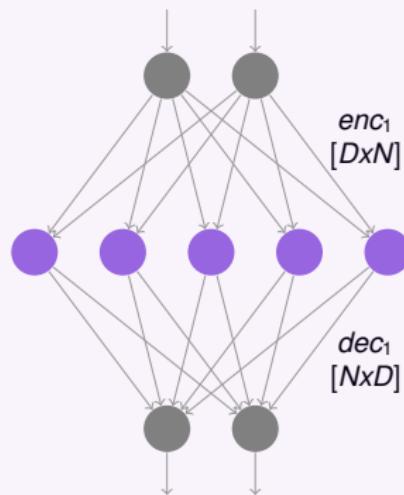
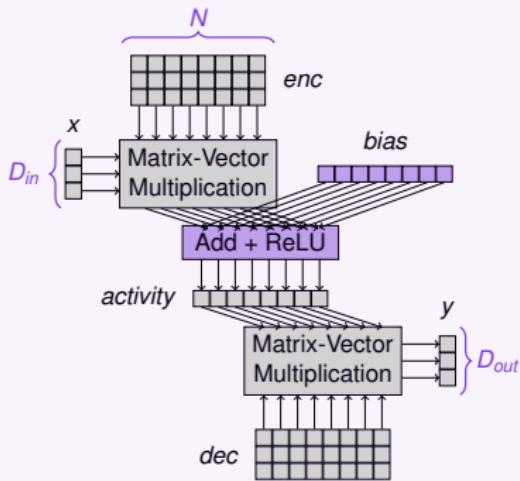


# Context

## Neural Engineering Framework

This boils down to two matrix multiplies

- ▶ One for encode
  - ↳ Plus neuron model (ReLU)

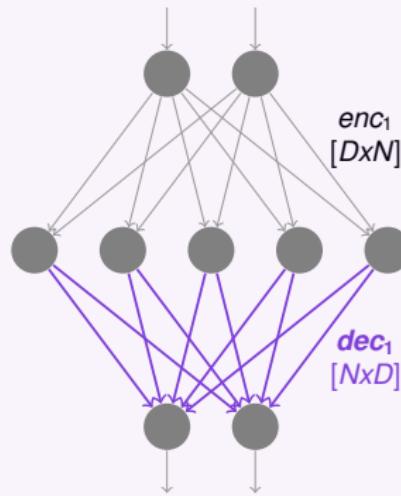
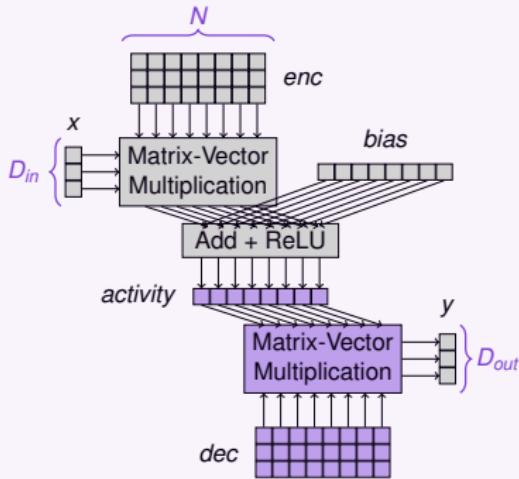


# Context

## Neural Engineering Framework

This boils down to two matrix multiplies

- ▶ One for encode
  - ↳ Plus neuron model (ReLU)
- ▶ One for decode

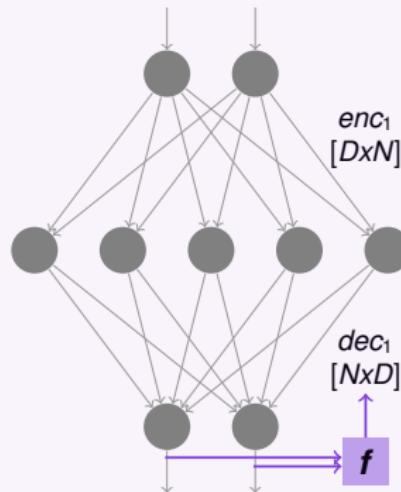
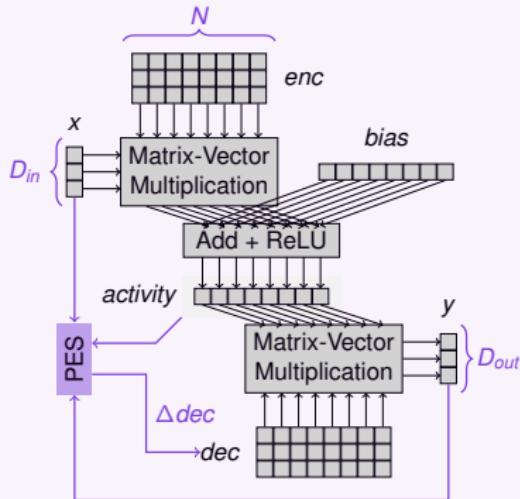


# Context

## Neural Engineering Framework

This boils down to two matrix multiplies

- ▶ One for encode
  - ↳ Plus neuron model (ReLU)
- ▶ One for decode
  - ↳ Plus **online learning**



## Motivation

---

GPU is great for matrix multiplies!

- ▶ Jetson TX1 can do 1 TFLOP/s

## Motivation

---

GPU is great for matrix multiplies!

- ▶ Jetson TX1 can do 1 TFLOP/s

However...

- ▶ No batching
- ▶ Limited on-chip memory
- ▶ No direct I/O access
- ▶ Little flexibility in data types

## Motivation

---

In Addition,

```
1 for i < N do
2   for j <  $D_{in}$  do
3     |  $a_i += x_j * enc_{ij}$ 
4   end
5 end
```

A single matrix multiply is simple to parallelize...

## Motivation

---

In Addition,

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for j < Dout do
7   for i < N do
8     | yj += ai * decji
9   end
10 end
```

A single matrix multiply is simple to parallelize...

...but two multiplies with different structure is less obvious.

## — Implementation

---

### Challenges

The challenges we address are:

- ▶ Parallelization with HLS
  - ↳ Temporal dependency (learning) — **no batching!**
  - ↳ Structural difference between encode & decode
- ▶ Resource limitations on-chip

## — Implementation

---

### Challenges

The challenges we address are:

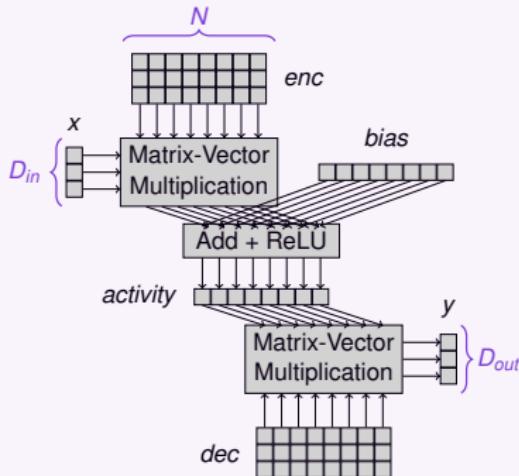
- ▶ Parallelization with HLS
  - ↳ Temporal dependency (learning) — **no batching!**
  - ↳ Structural difference between encode & decode
- ▶ Resource limitations on-chip

# Implementation

## High-Level Synthesis Parallelization

First let's estimate required cycles.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for j < Dout do
7   for i < N do
8     | yj += ai * decji
9   end
10 end
```

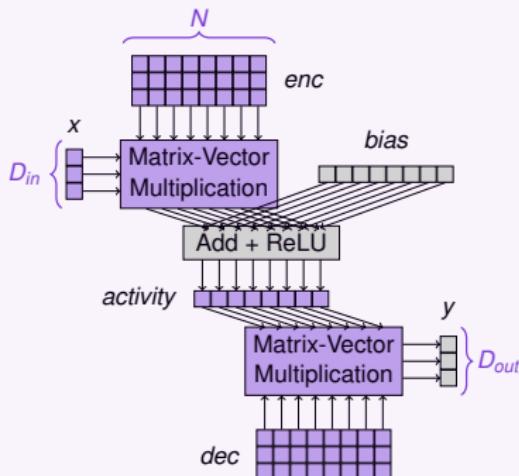


## Implementation

### High-Level Synthesis Parallelization

First let's estimate required cycles.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for j < Dout do
7   for i < N do
8     | yj += ai * decji
9   end
10 end
```



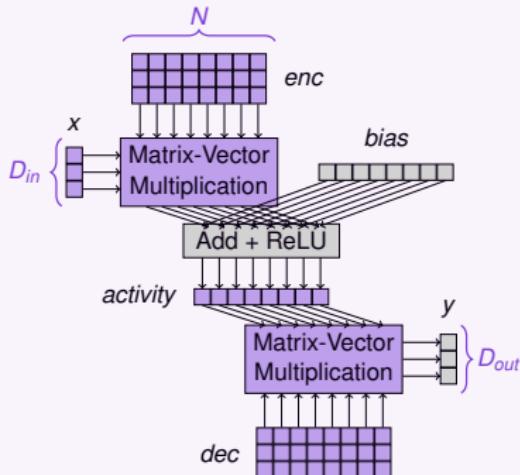
- ▶  $2 \times N \times D$

# Implementation

## High-Level Synthesis Parallelization

First let's estimate required cycles.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for j < Dout do
7   for i < N do
8     | yj += ai * decji
9   end
10 end
```



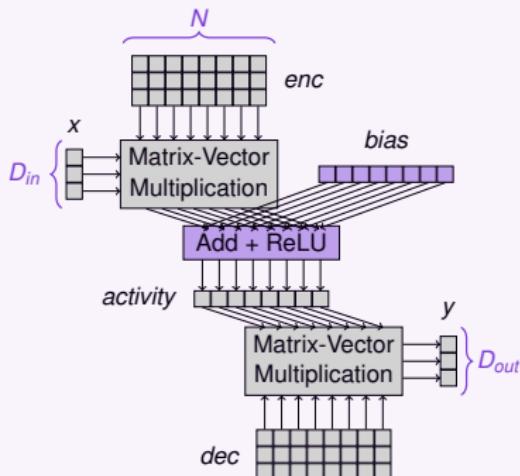
- ▶  $3 \times (2 \times N \times D)$

## Implementation

### High-Level Synthesis Parallelization

First let's estimate required cycles.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for j < Dout do
7   for i < N do
8     | yj += ai * decji
9   end
10 end
```



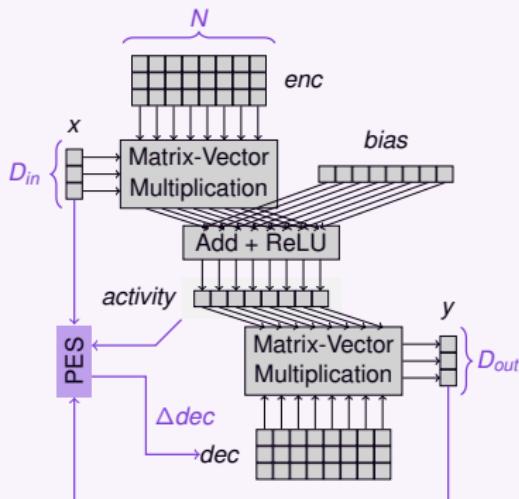
- ▶  $3 \times (2 \times N \times D) + 2 \times N$

# Implementation

## High-Level Synthesis Parallelization

First let's estimate required cycles.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for j < Dout do
7   for i < N do
8     | yj += ai * decji
9   end
10 end
```



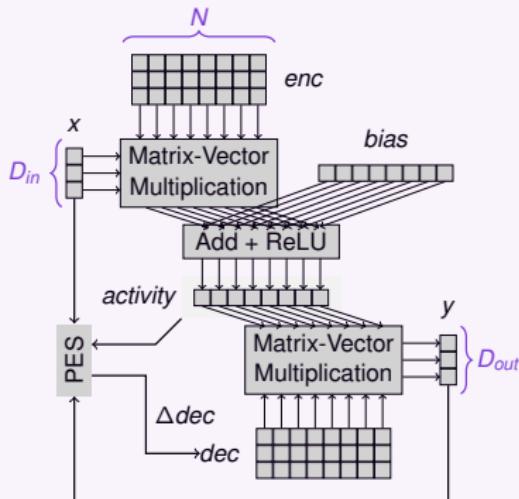
$$\triangleright 3 \times (3 \times N \times D) + 2 \times N$$

# Implementation

## High-Level Synthesis Parallelization

First let's estimate required cycles.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for j < Dout do
7   for i < N do
8     | yj += ai * decji
9   end
10 end
```



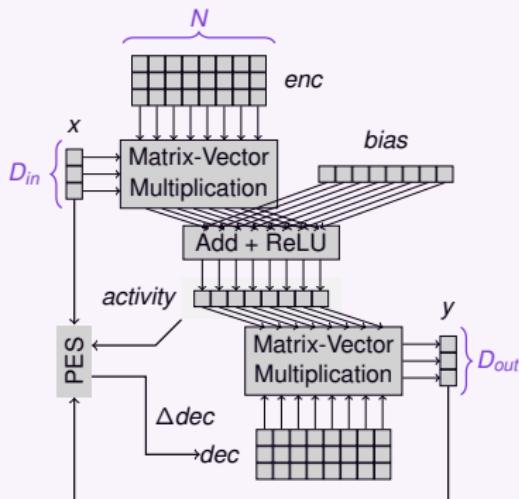
- ▶  $3 \times (3 \times N \times D) + 2 \times N$
- ▶ e.g.  $N=200$  and  $D=2$   
↳ = 4000 cycles

# Implementation

## High-Level Synthesis Parallelization

First let's estimate required cycles.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for j < Dout do
7   for i < N do
8     | yj += ai * decji
9   end
10 end
```



- ▶  $3 \times (3 \times N \times D) + 2 \times N$
- ▶ e.g.  $N=200$  and  $D=2$ 
  - ↪ = 4000 cycles

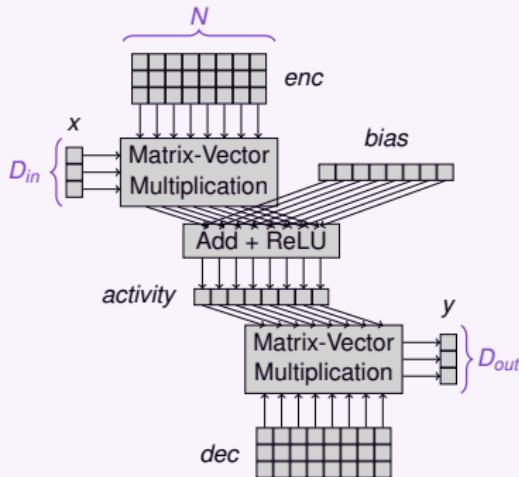
but HLS estimates 6000 cycles with naive approach!

## Implementation

### High-Level Synthesis Parallelization

Now we restructure the second matrix multiply.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for j < Dout do
7   for i < N do
8     | yj += ai * decji
9   end
10 end
```

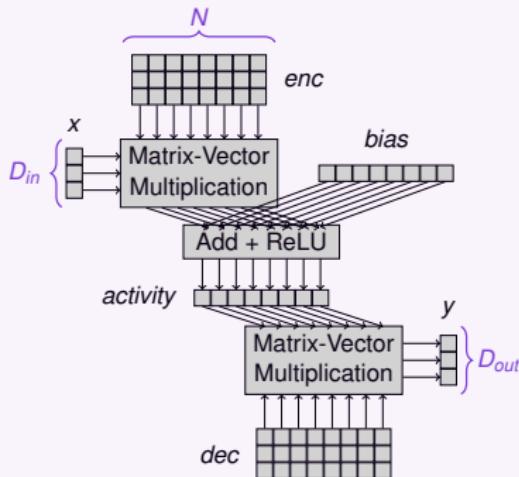


## Implementation

### High-Level Synthesis Parallelization

Now we restructure the second matrix multiply.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for i < N do
7   for j < Dout do
8     | yj += ai * decji
9   end
10 end
```

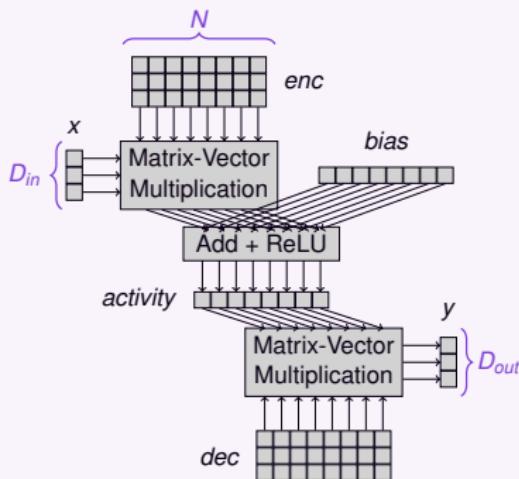


## Implementation

### High-Level Synthesis Parallelization

Now we restructure the second matrix multiply.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for i < N do
7   for j < Dout do
8     | yj += ai * decji
9   end
10 end
```



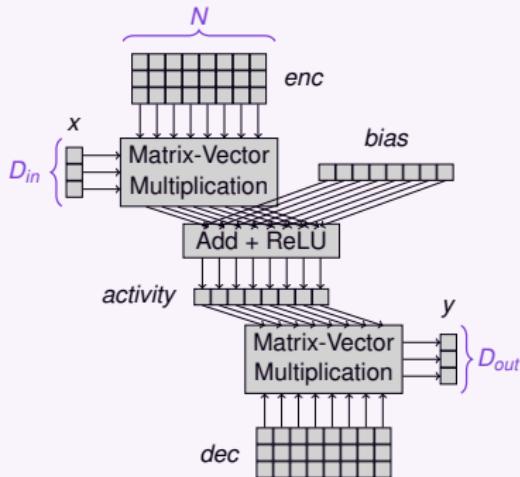
This improved II and got us to 4800 cycles.  
(recall we estimated 4000)

## Implementation

### High-Level Synthesis Parallelization

Now we restructure the second matrix multiply.

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for i < N do
7   for j < Dout do
8     | yj += ai * decji
9   end
10 end
```



This improved II and got us to 4800 cycles.  
(recall we estimated 4000)

## — Implementation

---

### High-Level Synthesis Parallelization

Now loop structure is harmonized, we want to parallelize

## Implementation

---

### High-Level Synthesis Parallelization

Now loop structure is harmonized, we want to parallelize

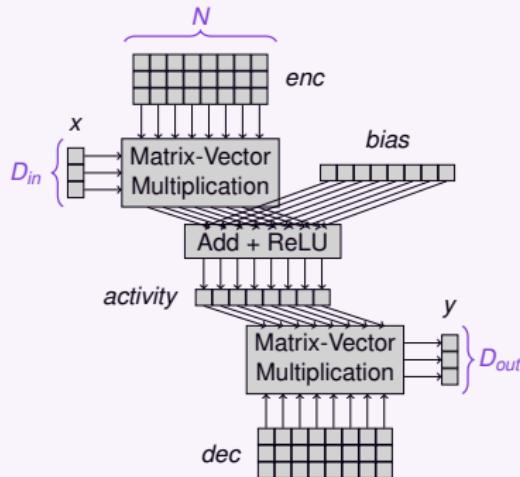
But `unroll` and `dataflow` pragmas could not automatically extract parallelism

# Implementation

## High-Level Synthesis Parallelization

We add explicit parallel structure

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
6 for i < N do
7   for j < Dout do
8     | yj += ai * decji
9   end
10 end
```

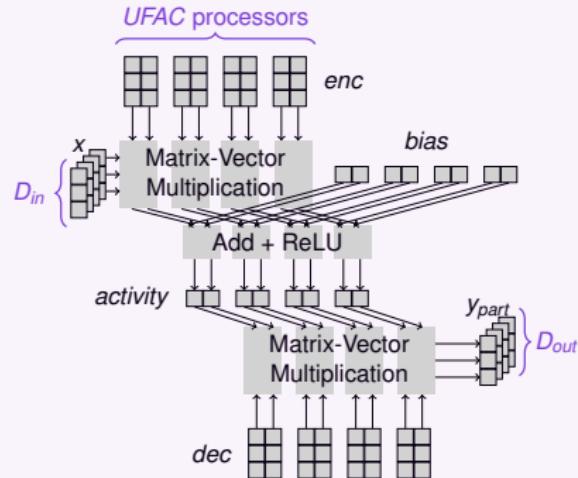


# Implementation

## High-Level Synthesis Parallelization

We add explicit parallel structure

```
1 for  $k < UFAC$  do
2   for  $i < \lceil \frac{N}{UFAC} \rceil$  do
3     for  $j < D_{in}$  do
4        $a_i += x_{kj} * enc_{ij}$ 
5     end
6   end
7   for  $i < \lceil \frac{N}{UFAC} \rceil$  do
8     for  $j < D_{out}$  do
9        $y_{kj}^{part} += a_i * dec_{ji}$ 
10    end
11  end
12 end
```

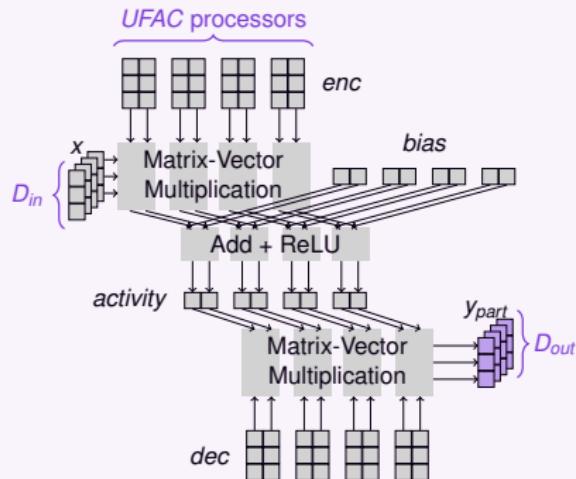


# Implementation

## High-Level Synthesis Parallelization

Which creates partial results that must be accumulated

```
1  for k < UFAC do
2      for i < ⌈ N / UFAC ⌉ do
3          for j < Din do
4              | ai += xkj * encij
5          end
6      end
7      for i < ⌈ N / UFAC ⌉ do
8          for j < Dout do
9              | ykjpart += ai * decji
10         end
11     end
12 end
```

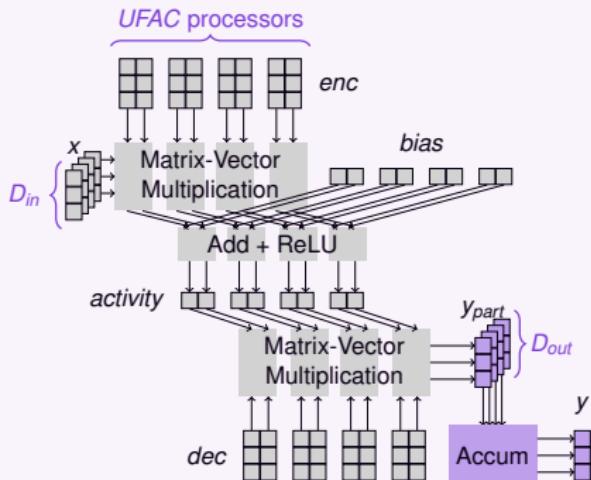


# Implementation

## High-Level Synthesis Parallelization

Which creates partial results that must be accumulated

```
1  for k < UFAC do
2      for i < ⌈ N / UFAC ⌉ do
3          for j < Din do
4              | ai += xkj * encij
5          end
6      end
7      for i < ⌈ N / UFAC ⌉ do
8          for j < Dout do
9              | ykjpart += ai * decji
10         end
11     end
12   end
13   for j < Dout do
14       for k < UFAC do
15           | yj += ykjpart
16       end
17   end
```

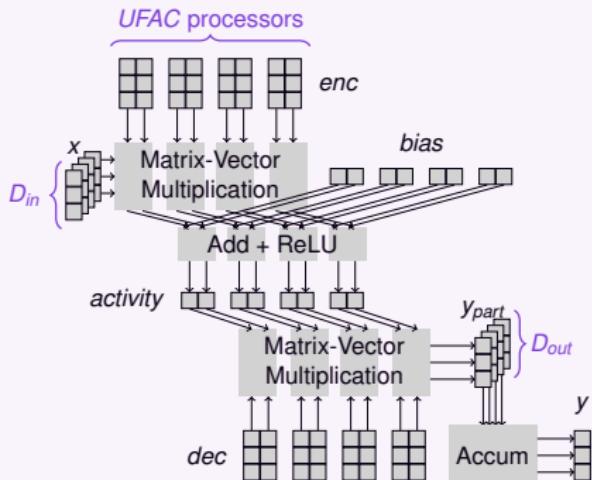


# Implementation

## High-Level Synthesis Parallelization

Which creates partial results that must be accumulated

```
1 for k < UFAC do
2   for i < ⌈ N / UFAC ⌉ do
3     for j < Din do
4       | ai += xkj * encij
5     end
6   end
7   for i < ⌈ N / UFAC ⌉ do
8     for j < Dout do
9       | ykjpart += ai * decji
10    end
11  end
12 end
13 for j < Dout do
14   for k < UFAC do
15     | yj += ykjpart
16   end
17 end
```

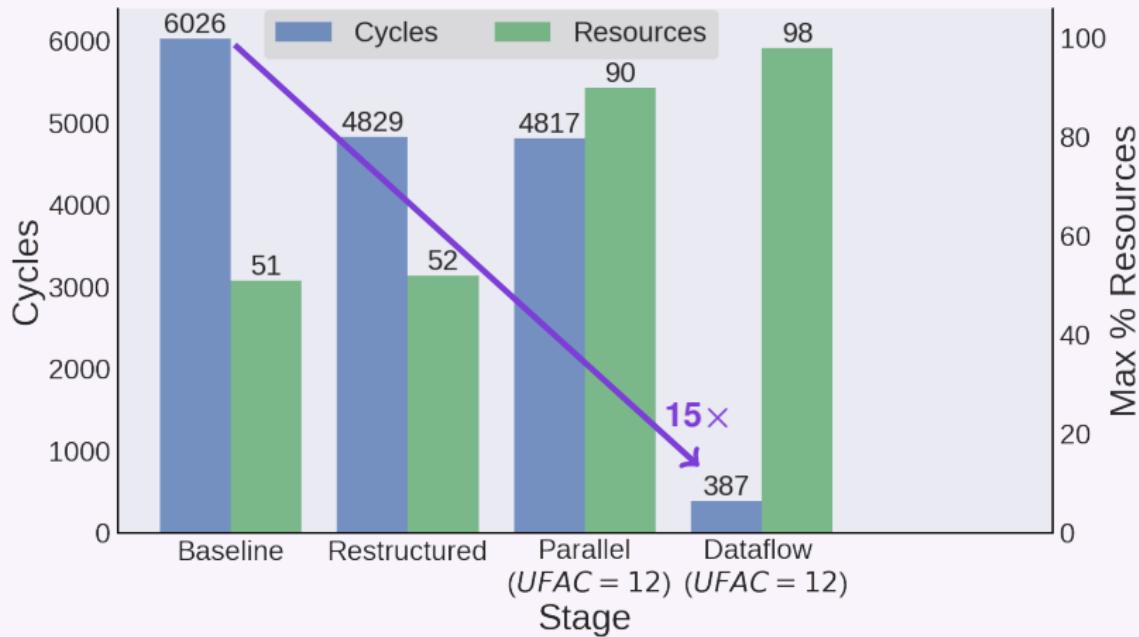


400 cycles with  $UFAC = 12$ .

## Results

### HLS Progression

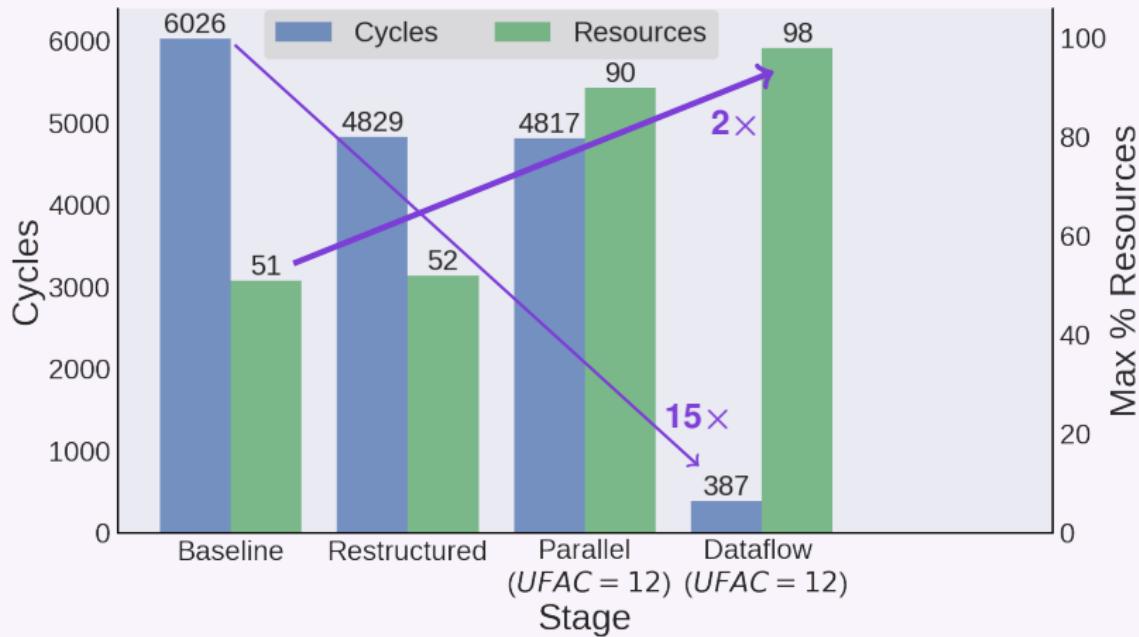
Clever HLS structure allows us to use the entire chip



## Results

### HLS Progression

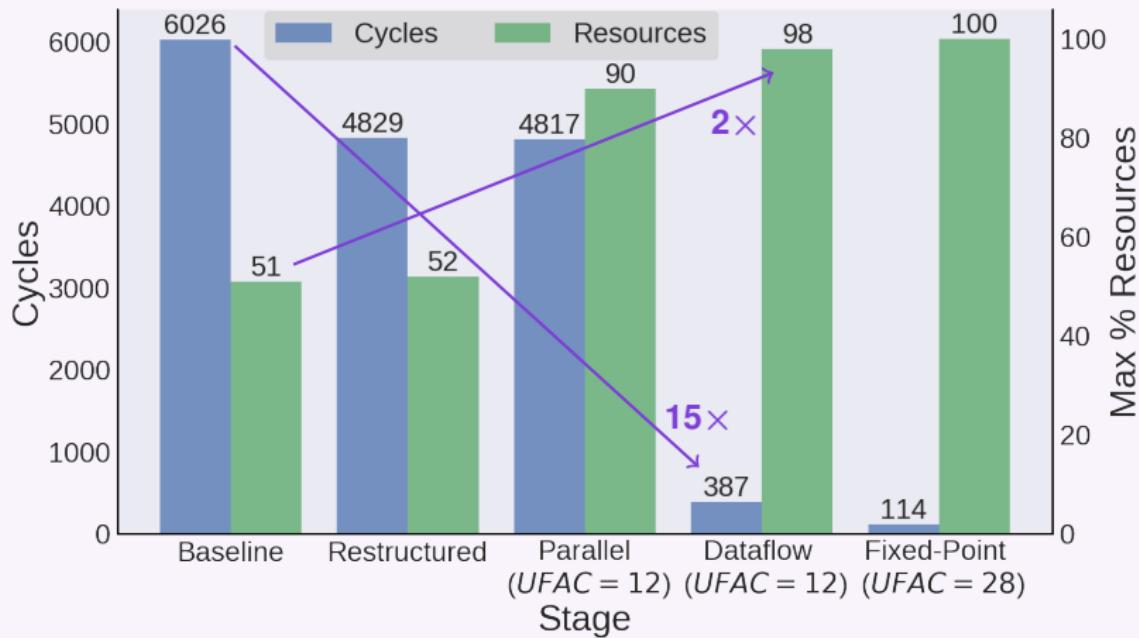
Clever HLS structure allows us to use the entire chip



# Results

## HLS Progression

Fixed-point further shrinks compute tile



## — Implementation

---

### Challenges

The challenges we address are:

- ▶ Parallelization with HLS
  - ↳ Temporal dependency (learning) — **no batching!**
  - ↳ Structural difference between encode & decode
- ▶ Resource limitations on-chip

## Implementation

---

### Fixed-Point Tuning

We create 5 different fixed-point types using the `ap_fixed` template type provided by Vivado HLS

## Implementation

---

### Fixed-Point Tuning

We create 5 different fixed-point types using the `ap_fixed` template type provided by Vivado HLS

For each type we have many choices for:

- ▶ Number of word bits
- ▶ Number of integer bits
- ▶ Rounding strategy
- ▶ Overflow strategy

## Implementation

---

### Fixed-Point Tuning

We create 5 different fixed-point types using the `ap_fixed` template type provided by Vivado HLS

For each type we have many choices for:

- ▶ Number of word bits
- ▶ Number of integer bits
- ▶ Rounding strategy
- ▶ Overflow strategy

This results in  $\approx 10^{20}$  configurations!

## — Fixed-Point

---

Hyperopt

We use a Python package called [Hyperopt](#) to automatically tune our hyper-parameters and extract Pareto optimal designs for a given cost function.

## — Fixed-Point

---

Hyperopt

We want to minimize error and maximize performance

## — Fixed-Point

---

Hyperopt

We want to minimize error and maximize performance

We use Hyperopt to minimize **error** and **resources**

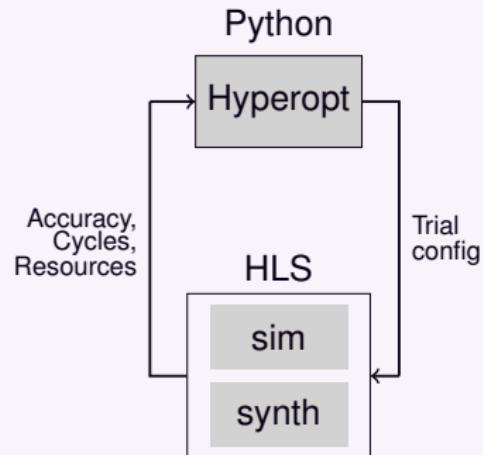
- ▶ Small, high-performance tile

## Fixed-Point

### Hyperopt

Each Hyperopt trial runs:

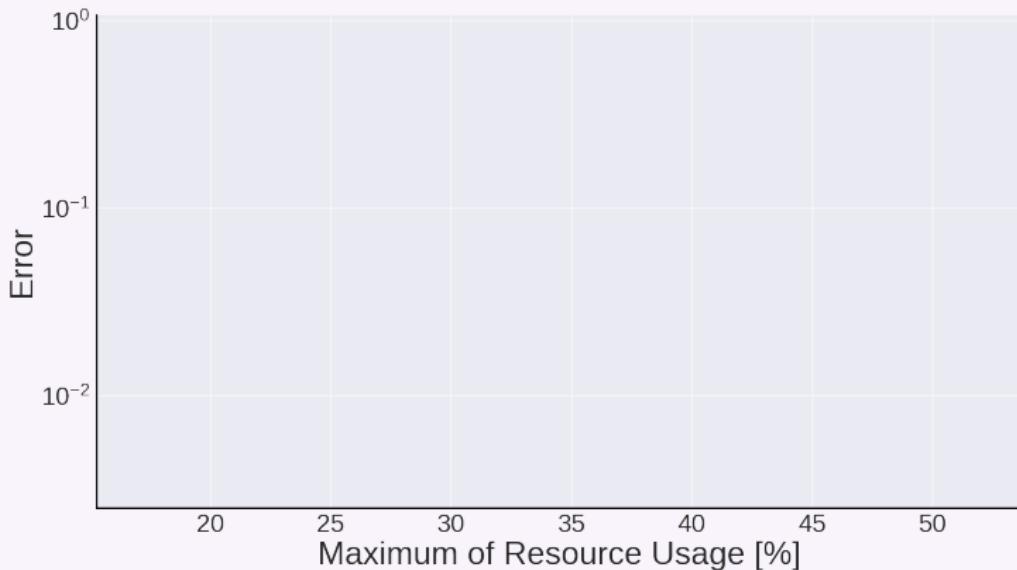
- ▶ C simulation
- ▶ C synthesis (no P & R)



# Results

## Hyperopt

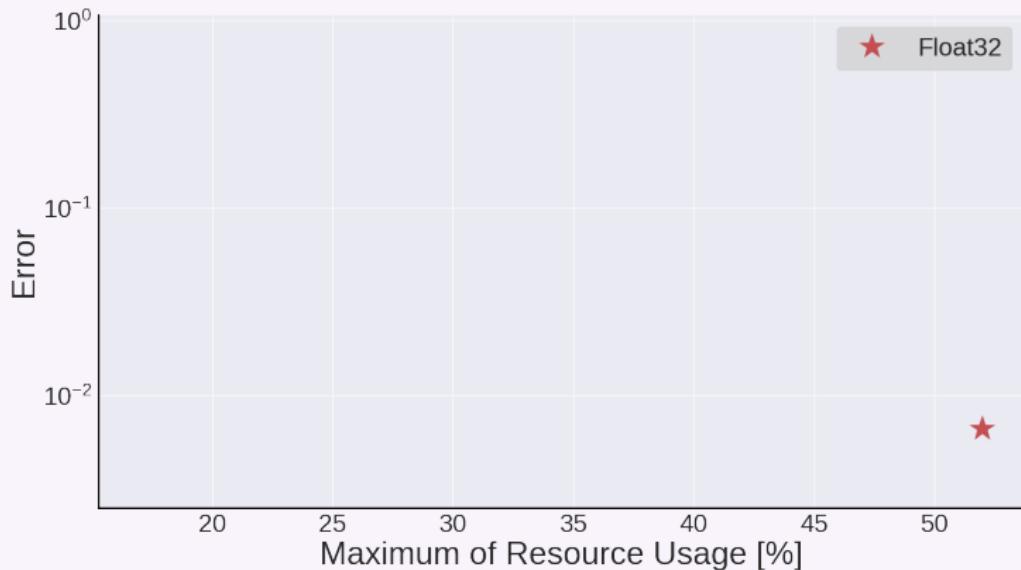
Let's look at some results:



# Results

## Hyperopt

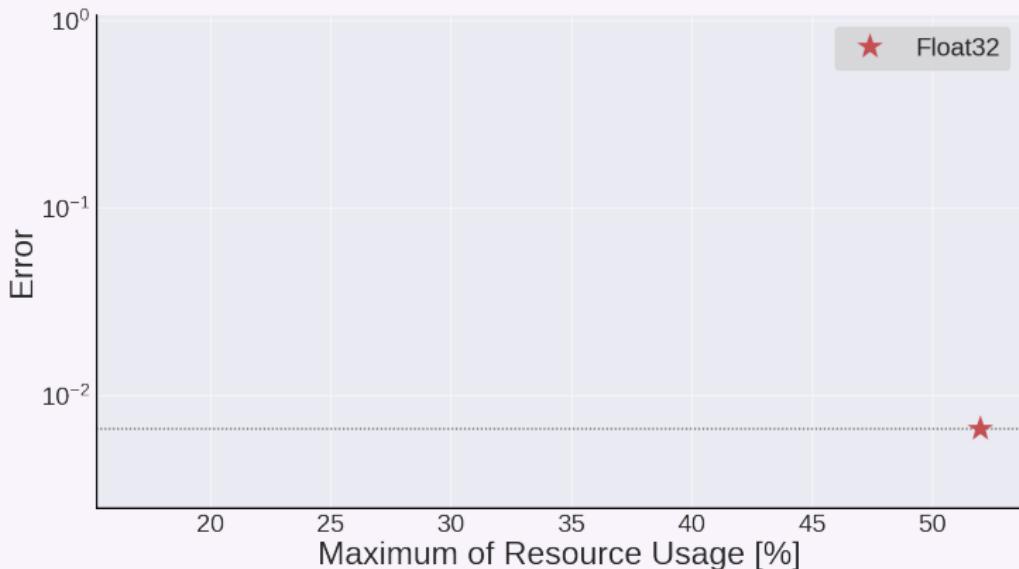
Let's look at some results:



# Results

## Hyperopt

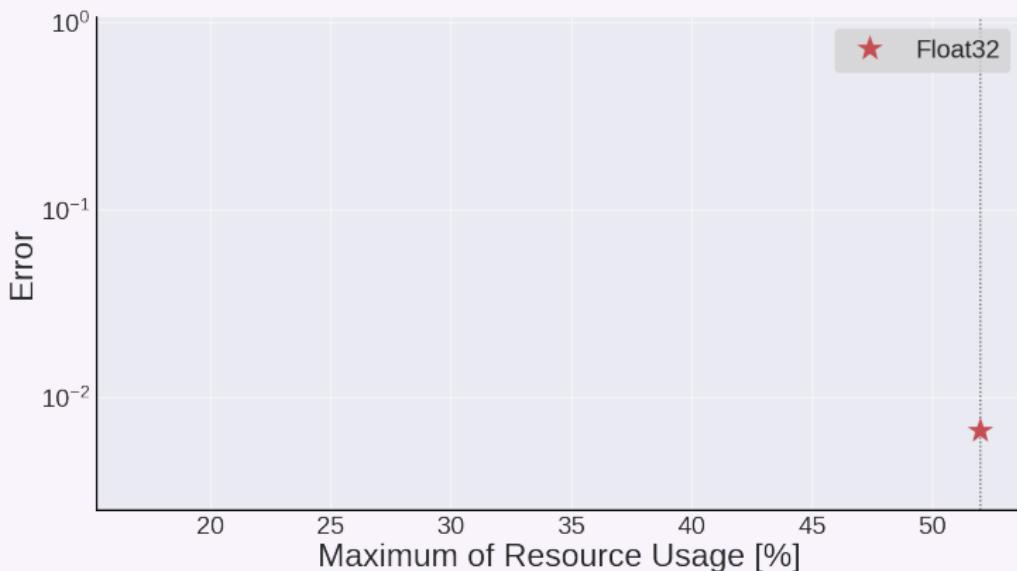
Let's look at some results:



# Results

## Hyperopt

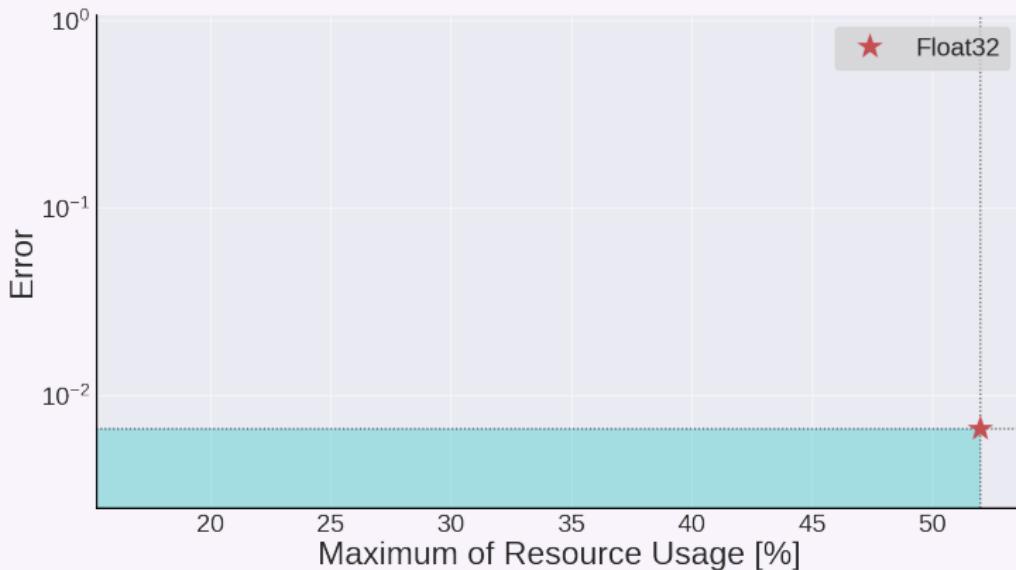
Let's look at some results:



# Results

## Hyperopt

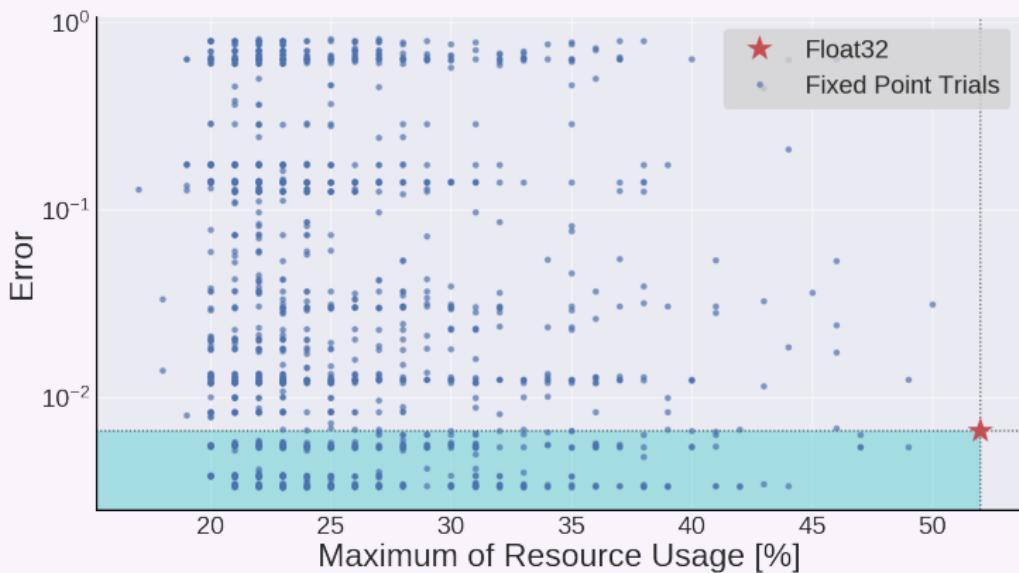
Let's look at some results:



# Results

## Hyperopt

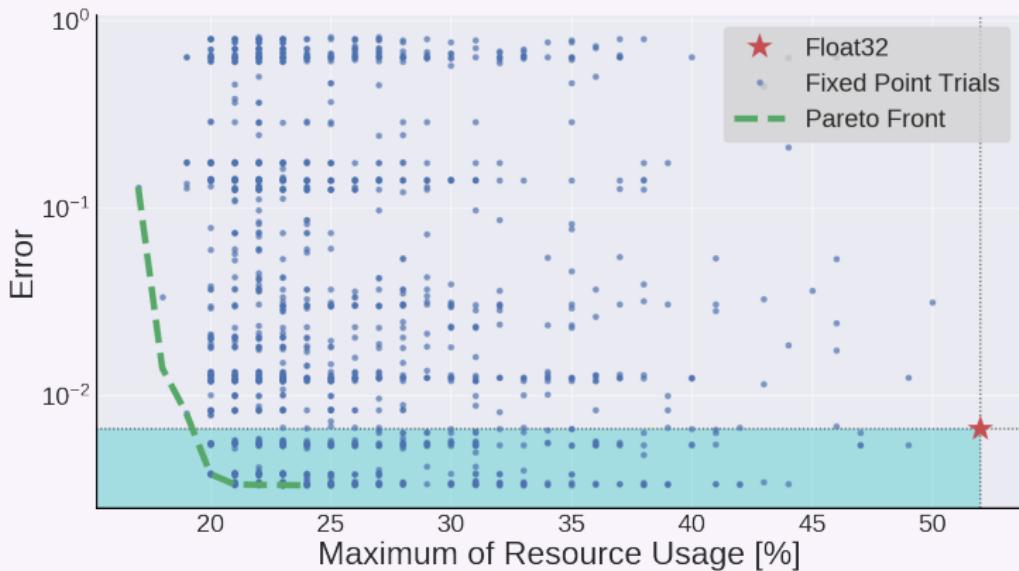
Let's look at some results:



# Results

## Hyperopt

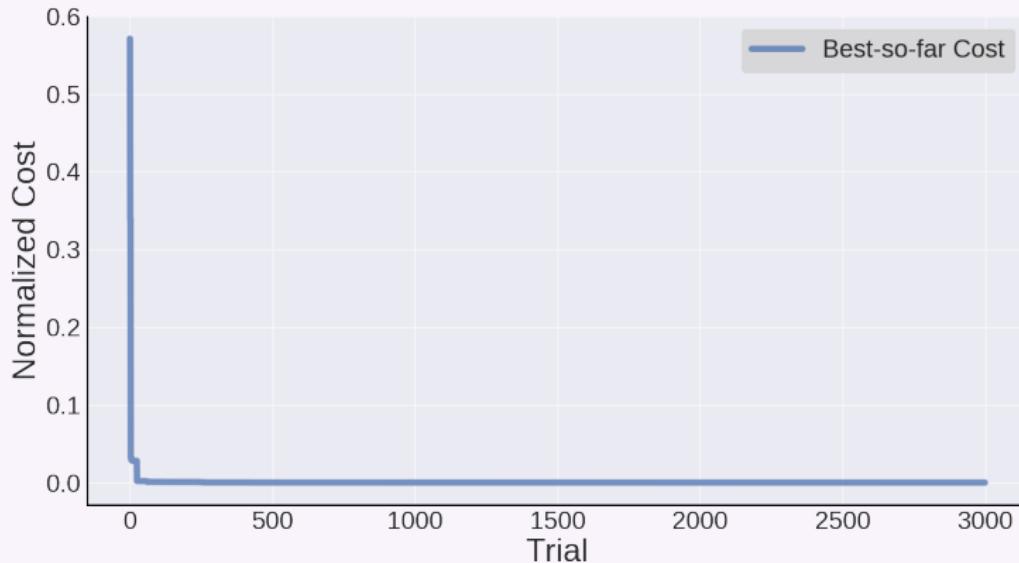
Let's look at some results:



## Results

### Hyperopt

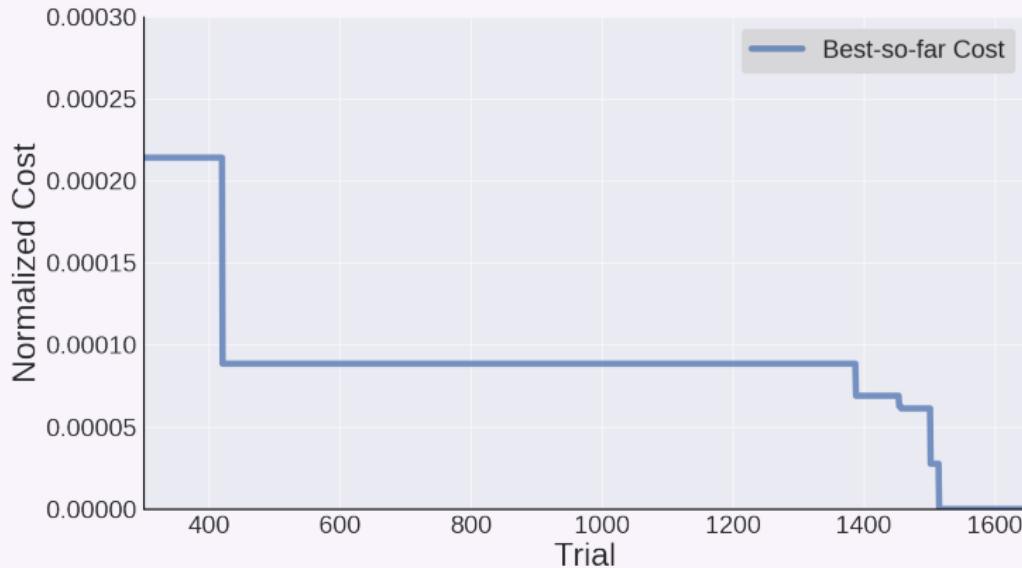
We let Hyperopt run for 3000 trials using a cost function that minimizes error and resources



## Results

### Hyperopt

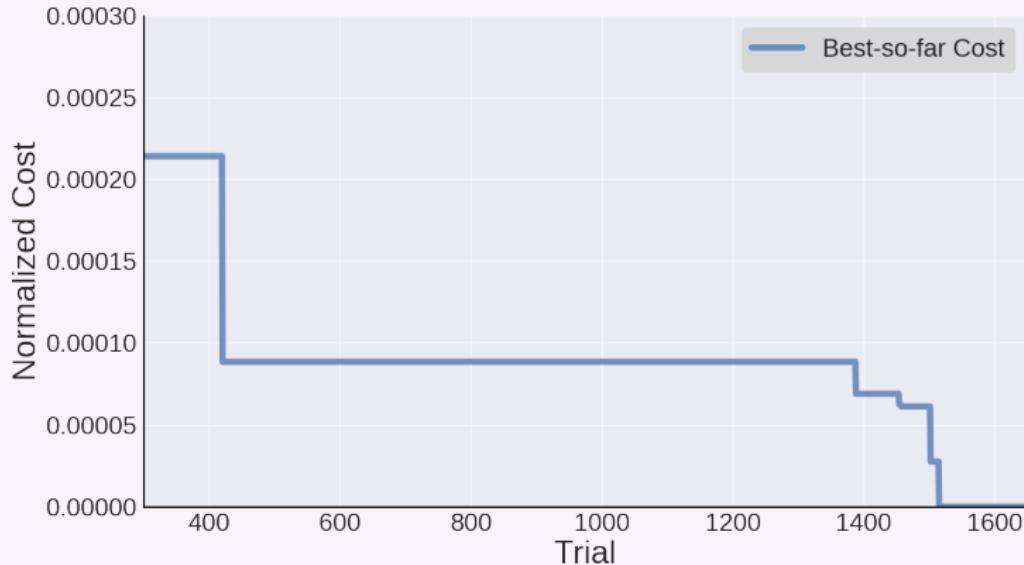
In fact, Hyperopt converges to the best solution after 1500 trials



## Results

### Hyperopt

And, we are within 99% after  $\approx 400$  trials which takes 4 hours  
(i7-6700k)

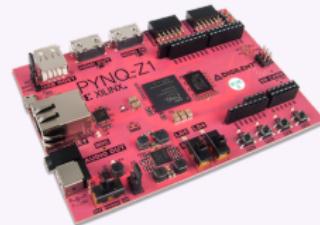


# Results

---

Performance vs. Jetson TX1

PYNQ uses a Zynq FPGA with 28nm technology.



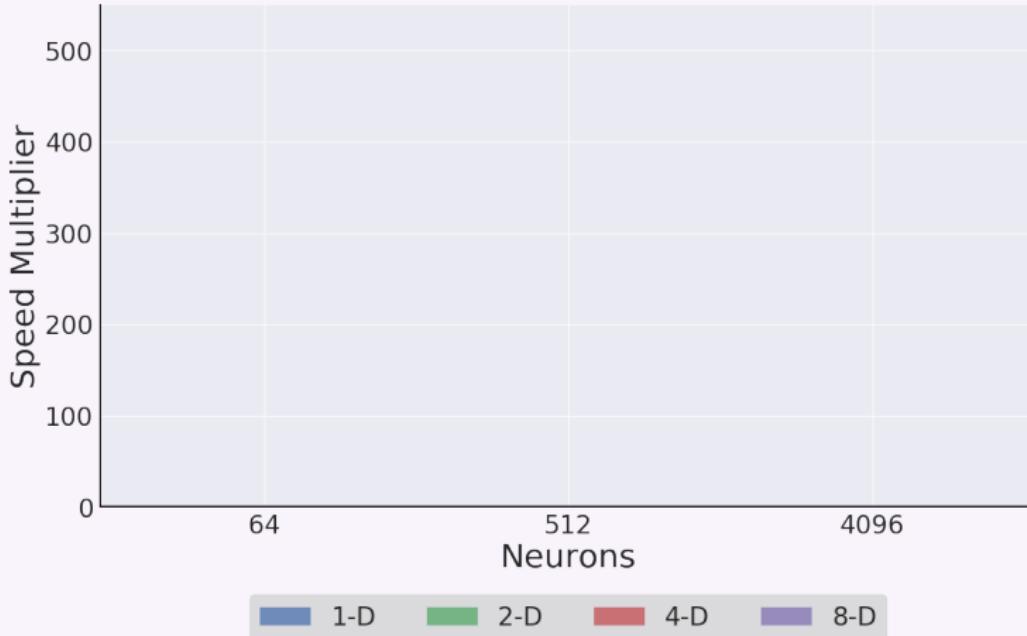
Jetson TX1 uses an embedded GPU with 20nm technology.



## Results

Performance vs. Jetson TX1

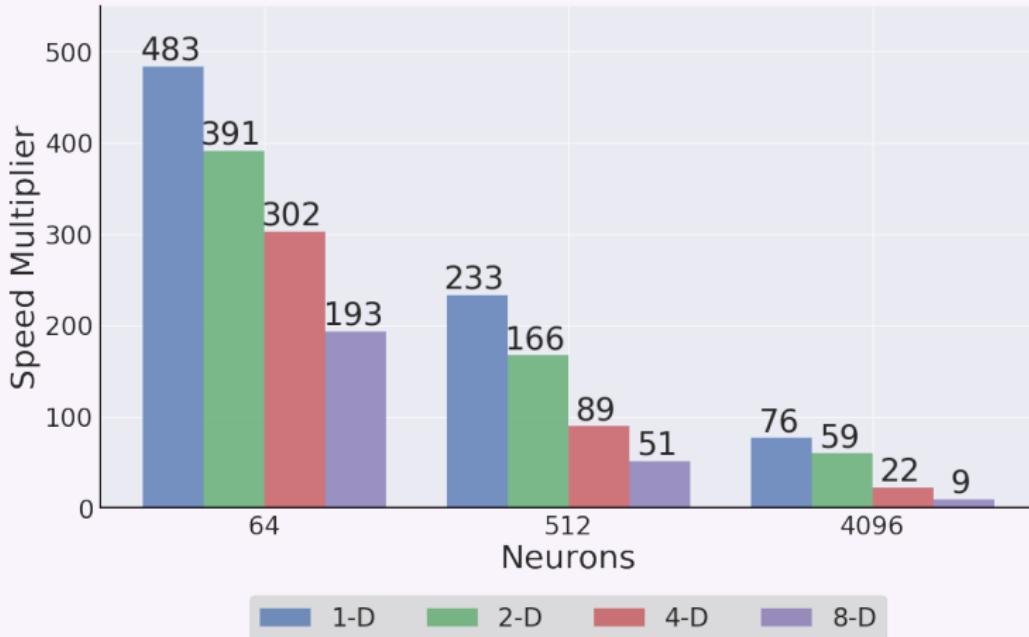
FPGA Speedup compared to cuBLAS implementation on Jetson TX1.



## Results

Performance vs. Jetson TX1

FPGA Speedup compared to cuBLAS implementation on Jetson TX1.



## — Practical Applications —

Our NengoFPGA PYNQ implementation can accomplish useful work

## Practical Applications

3000 neurons ( $D=6$ ) can form the basis for an adaptive motor controller [DeWolf et al., 2016]

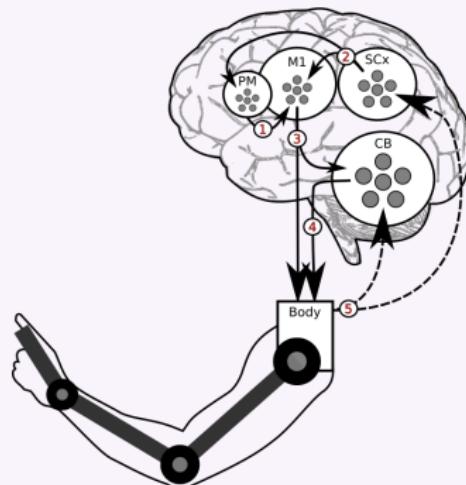
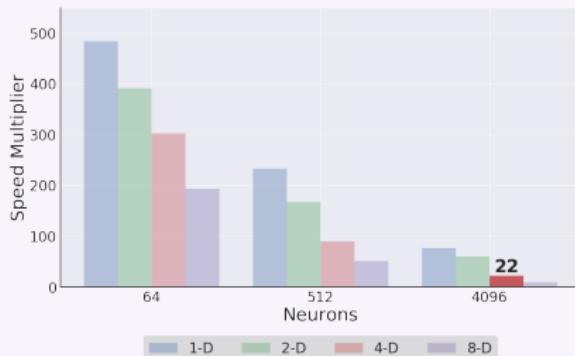


Image from "A spiking neural model of adaptive arm control" - DeWolf et al.

## Practical Applications

500 neurons ( $D=30$ ) can adapt and control a 15-joint body simulation [Stewart et al., 2015]

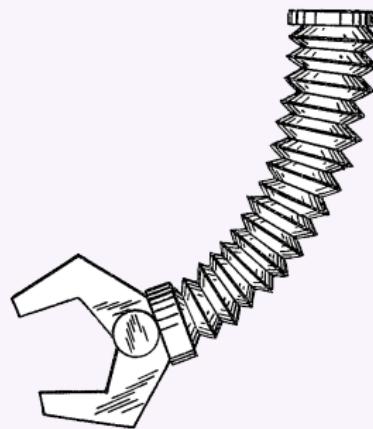
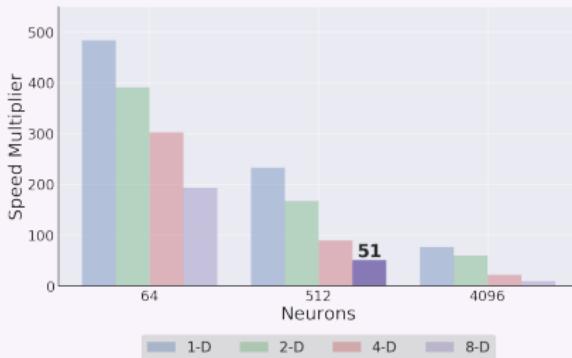


Image from Google

## Practical Applications

We control a variable mass inverted pendulum with comparable accuracy to the floating-point design ( $N=1000$ ,  $D=1$ )

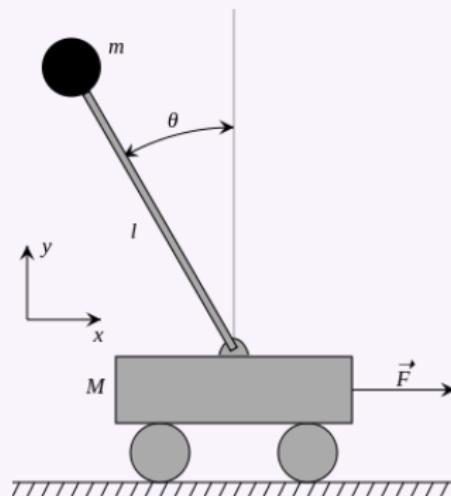
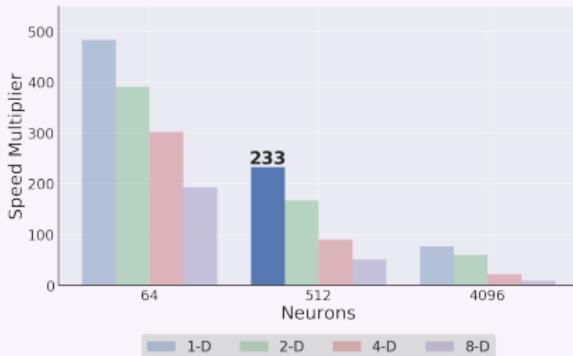


Image from Wikipedia

The core Nengo framework is directly integrated with the [NengoFPGA](#) Python package with two modes of operation:

- ▶ Direct embedded mode

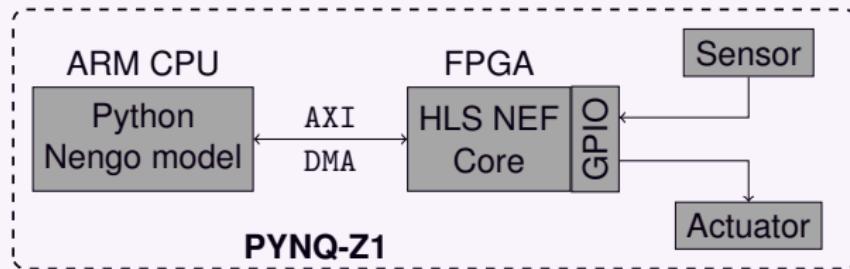
The core Nengo framework is directly integrated with the [NengoFPGA](#) Python package with two modes of operation:

- ▶ Direct embedded mode
- ▶ Remote PC control

## Interface

### Direct Mode

In direct mode we run Nengo models and drive the FPGA from the integrated ARM processor on the PYNQ SoC.



## — Interface

---

### Remote PC Mode

We also support running Nengo models on a host PC which provides:

- ▶ A familiar development environment

## — Interface

---

### Remote PC Mode

We also support running Nengo models on a host PC which provides:

- ▶ A familiar development environment
- ▶ Access to higher performance PC resources

## — Interface

---

### Remote PC Mode

We also support running Nengo models on a host PC which provides:

- ▶ A familiar development environment
- ▶ Access to higher performance PC resources
- ▶ Access to other subsystems connected to a PC

## — Interface

---

### Remote PC Mode

We also support running Nengo models on a host PC which provides:

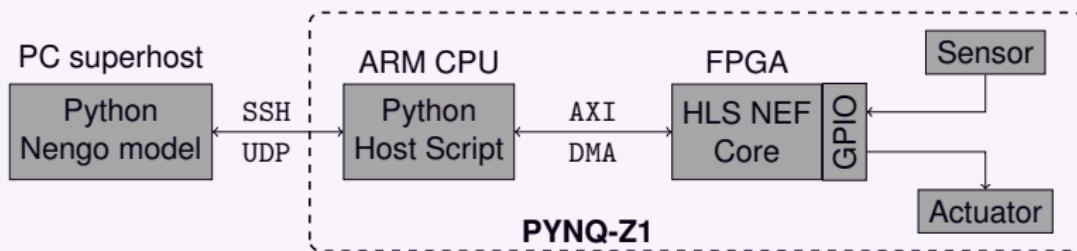
- ▶ A familiar development environment
- ▶ Access to higher performance PC resources
- ▶ Access to other subsystems connected to a PC

With only  $\approx 24\%$  communication overhead.

## Interface

### Remote PC Mode

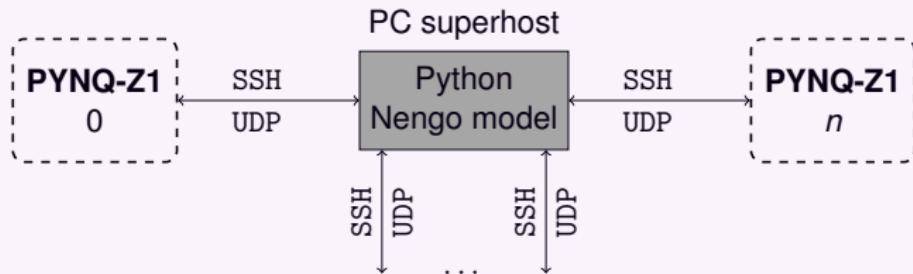
In this configuration, Nengo runs on the remote PC and the local ARM acts only to relay information



## Interface

### Remote PC Mode

Using a remote host also allows multiple devices to be integrated into a single system.



## Conclusion

---

NengoFPGA is a user-friendly, efficient Python-accessible implementation designed for dynamic NEF style neural networks and cognitive modelling that outperforms comparable devices.

## Conclusion

---

NengoFPGA is a user-friendly, efficient Python-accessible implementation designed for dynamic NEF style neural networks and cognitive modelling that outperforms comparable devices.

- ▶ Clever HLS architecture → 15× improvement

## Conclusion

---

NengoFPGA is a user-friendly, efficient Python-accessible implementation designed for dynamic NEF style neural networks and cognitive modelling that outperforms comparable devices.

- ▶ Clever HLS architecture →  $15\times$  improvement
- ▶ Automatically tuned precision →  $3\times$  improvement

## Conclusion

---

NengoFPGA is a user-friendly, efficient Python-accessible implementation designed for dynamic NEF style neural networks and cognitive modelling that outperforms comparable devices.

- ▶ Clever HLS architecture →  $15\times$  improvement
- ▶ Automatically tuned precision →  $3\times$  improvement
- ▶  $10\text{--}500\times$  faster than Jetson TX1

## Conclusion

---

NengoFPGA is a user-friendly, efficient Python-accessible implementation designed for dynamic NEF style neural networks and cognitive modelling that outperforms comparable devices.

- ▶ Clever HLS architecture →  $15\times$  improvement
- ▶ Automatically tuned precision →  $3\times$  improvement
- ▶  $10\text{--}500\times$  faster than Jetson TX1
- ▶  $2\text{--}10\times$  less power than Jetson TX1

## — Next Steps

---

We continue work to make **NengoFPGA** a seamless and fully featured backend for Nengo including:

- ▶ Larger FPGAs
- ▶ More Nengo support
- ▶ Dedicated accelerators
- ▶ Better CNN support

**NengoFPGA** is available from  
**Applied Brain Research**