

A detailed anatomical illustration of a human brain in lateral view. The amygdala is highlighted in pink, and the hippocampus is highlighted in yellow. Handwritten labels in the bottom left corner identify the 'Middle temporal gyrus' and 'Hippocampus'. The background shows the internal structures of the brain, including the cerebral cortex and white matter tracts.

# Implementing NEF Neural Networks on Embedded FPGAs

We develop a Python accessible PYNQ implementation that outperforms the Jetson TX1 GPU computing NEF style neural networks by **10–500×** using **2–10×** less power.

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

Our design:

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

- ▶ Machine learning is growing

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



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



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



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



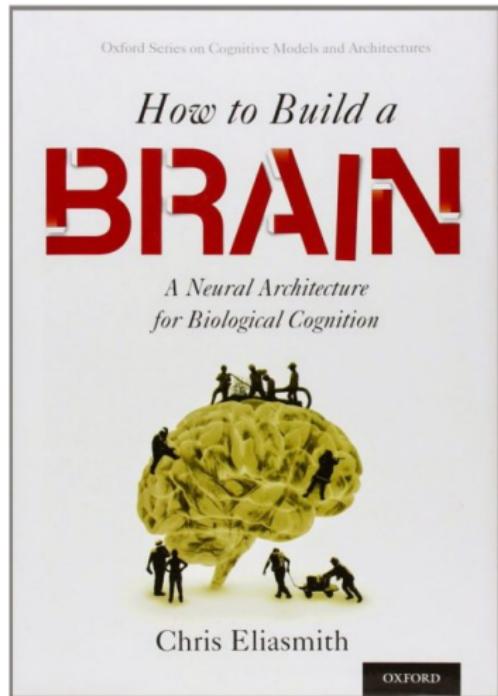
We target the [Neural Engineering Framework](#) (NEF) via the [Nengo](#) Python package

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

- ▶ Note NEF has different structure than DNN

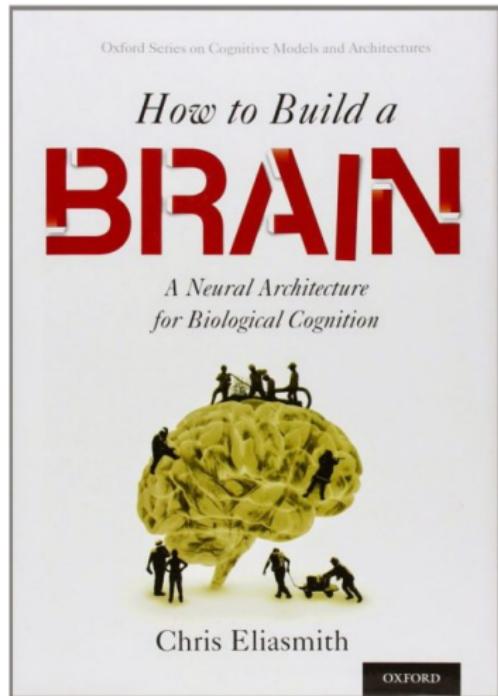
Nengo goes beyond DNN

- ▶ Dynamic models



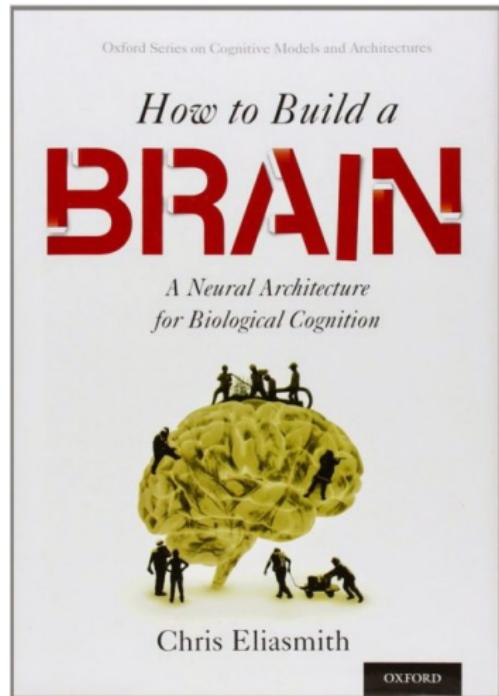
Nengo goes beyond DNN

- ▶ Dynamic models
- ▶ Real-time adaptive control



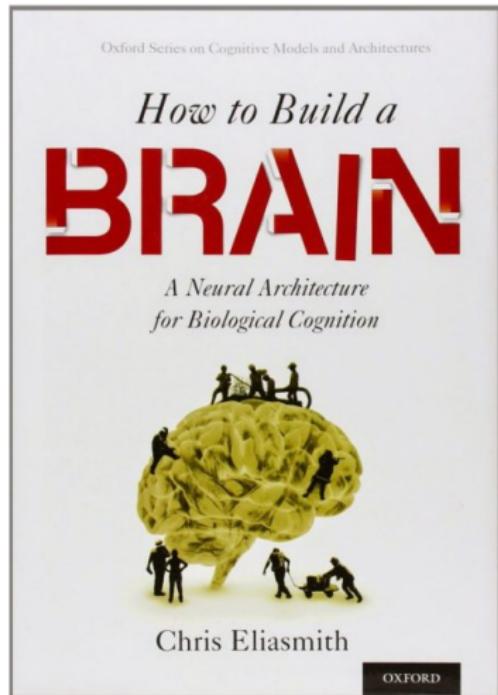
## Nengo goes beyond DNN

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

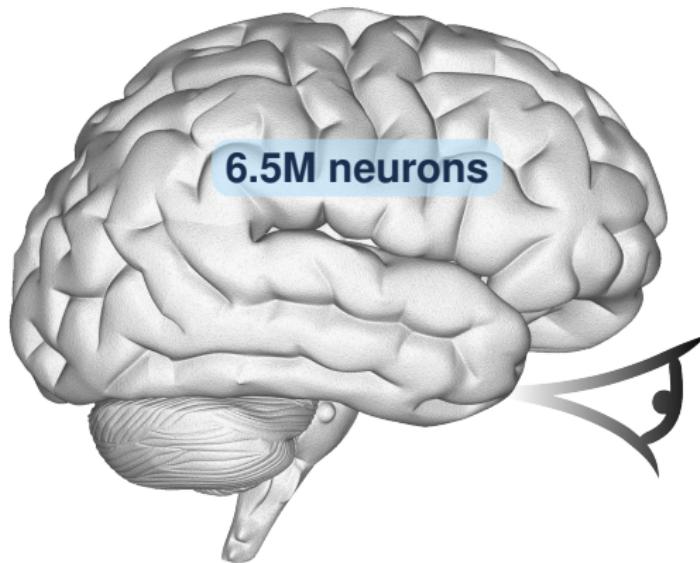


## Nengo goes beyond DNN

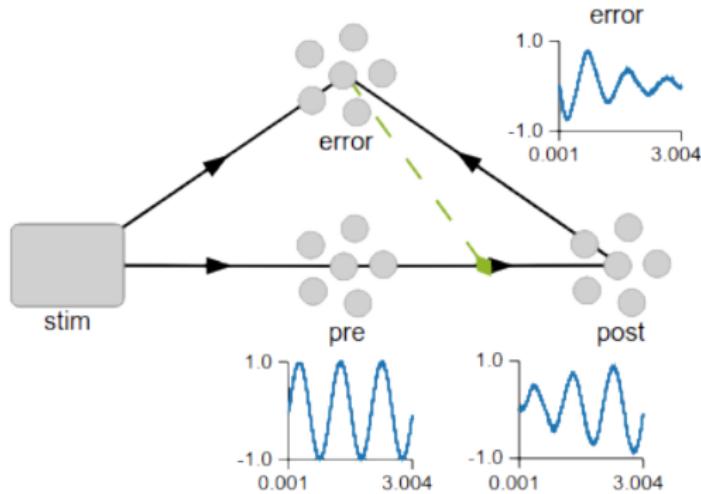
- ▶ Dynamic models
- ▶ Real-time adaptive control
- ▶ Biologically plausible cognitive architectures
- ▶ Hierarchical reinforcement learning

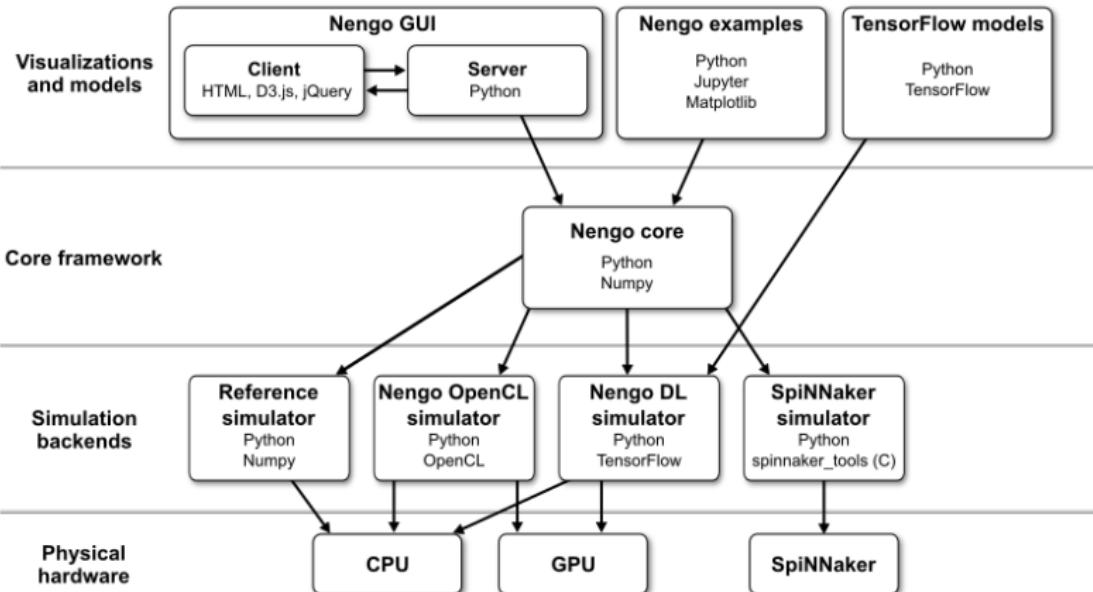


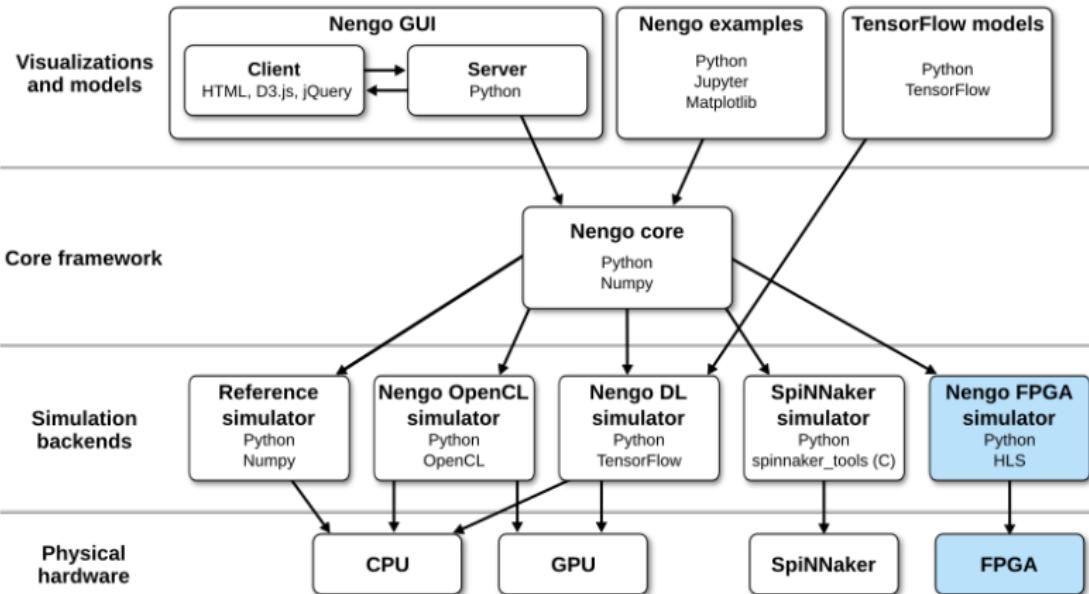
The world's largest functional brain model, Spaun, is built with the NEF and Nengo



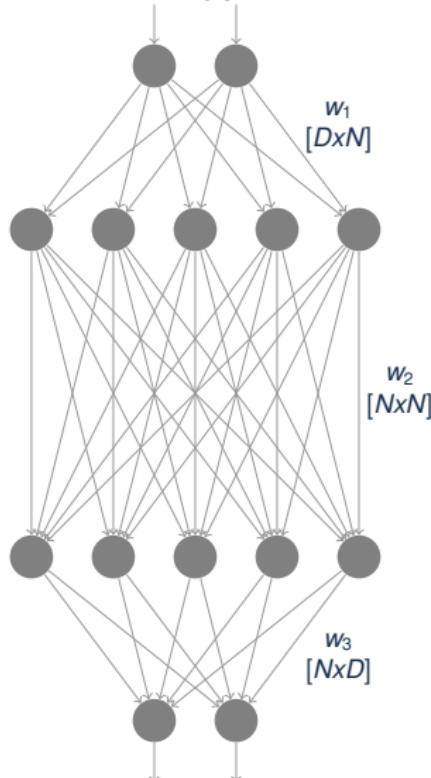
### Online learning displayed with Nengo GUI



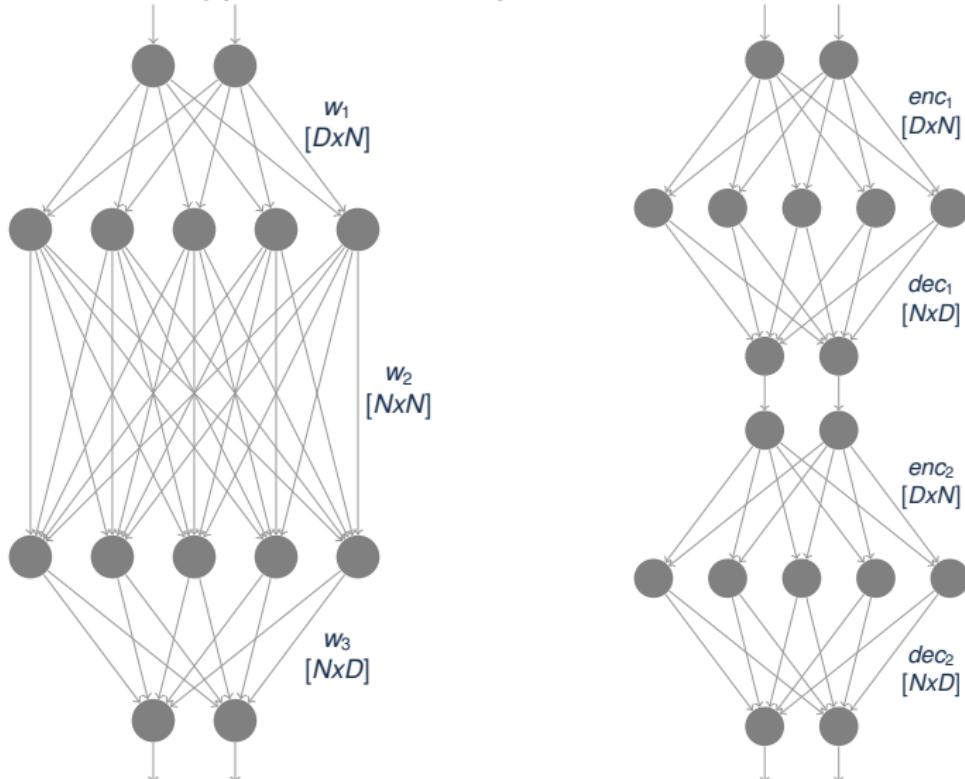




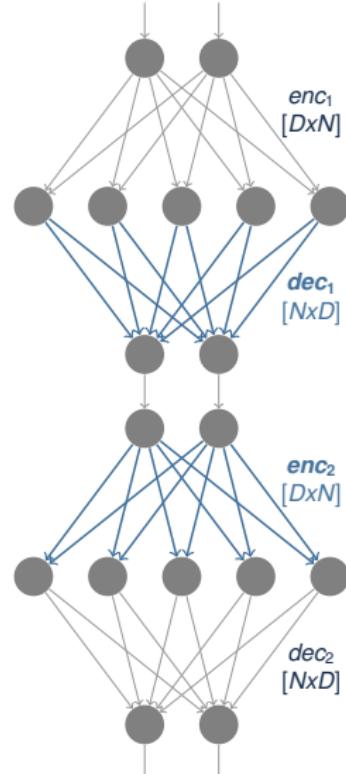
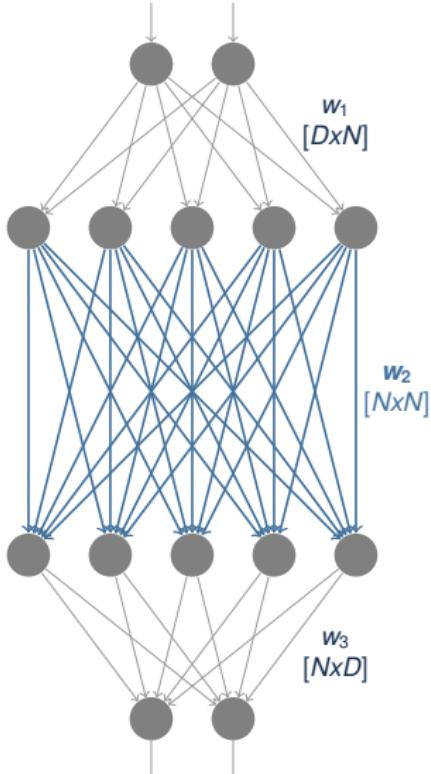
First let's look at typical DNN



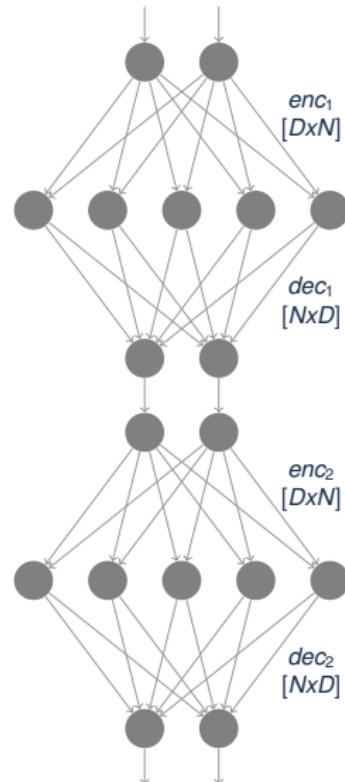
First let's look at typical DNN compared to an NEF network



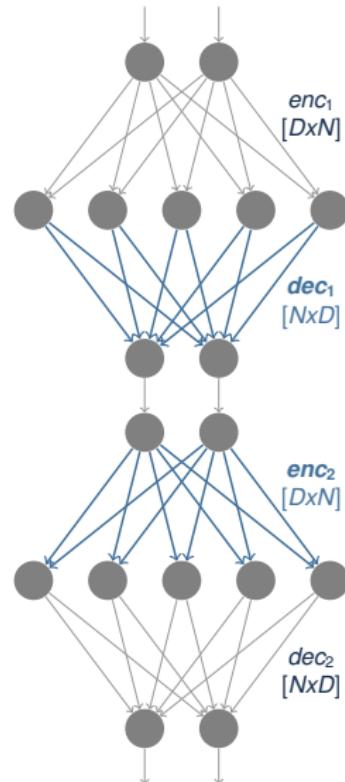
We factor the weight matrix into encode and decode matrices



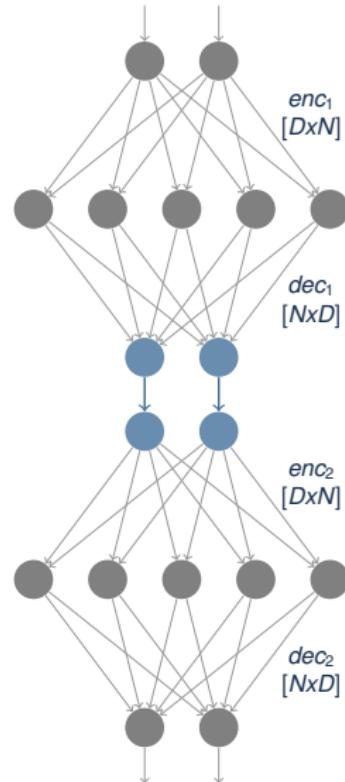
- We have  $N \gg D$



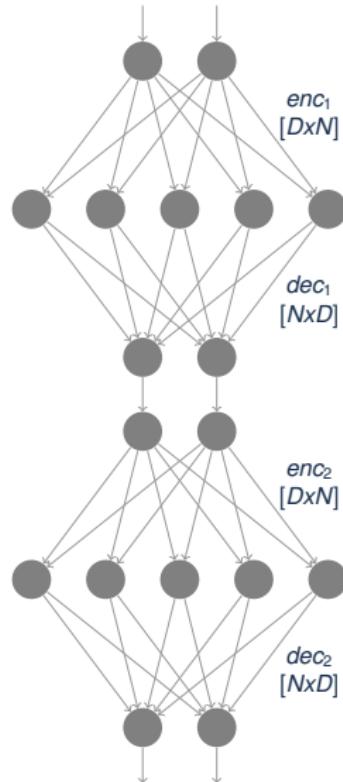
- We have  $N \gg D$ 
  - ↳ NEF saves **memory**



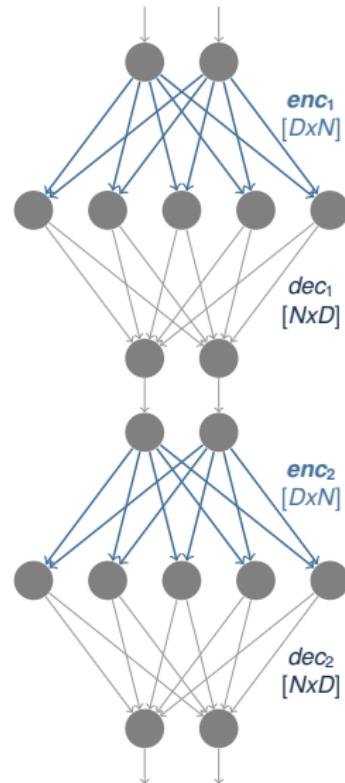
- ▶ We have  $N \gg D$ 
  - ↪ NEF saves memory
  - ↪ NEF saves bandwidth



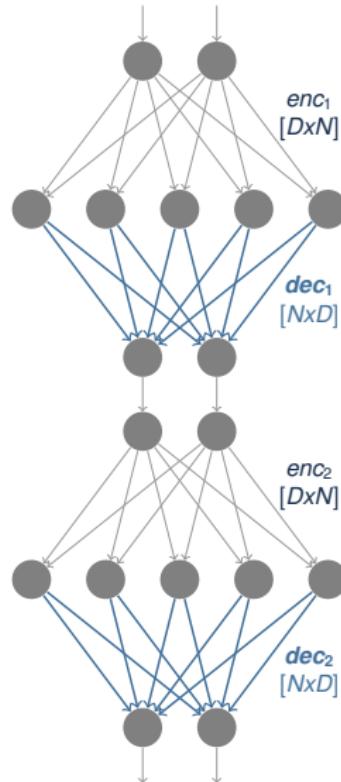
- ▶ We have  $N \gg D$ 
  - ↳ NEF saves memory
  - ↳ NEF saves bandwidth
- ▶ NEF adds temporal dimension



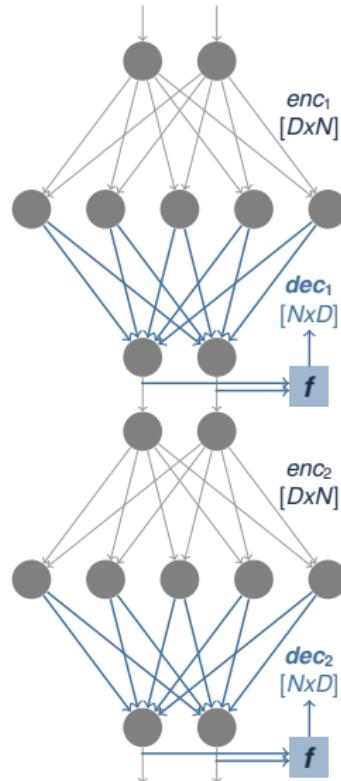
- ▶ Encoders are generated randomly



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

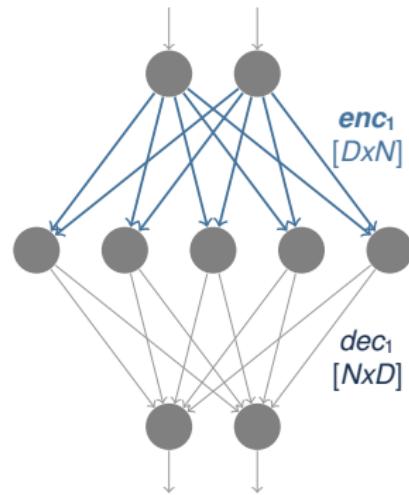
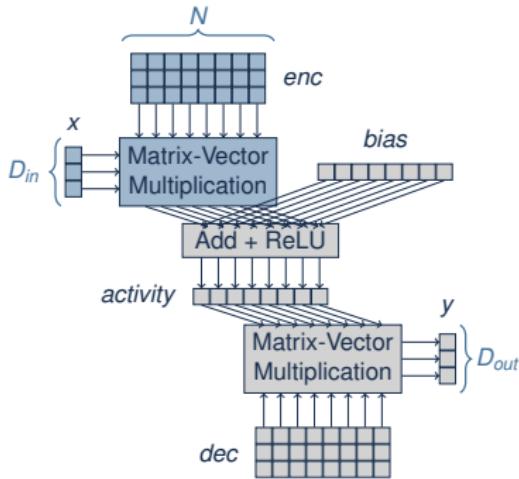


- ▶ Encoders are generated **randomly**
- ▶ Decoders can be:
  - ↪ Pre-solved in offline training and/or
  - ↪ Updated with **online learning**



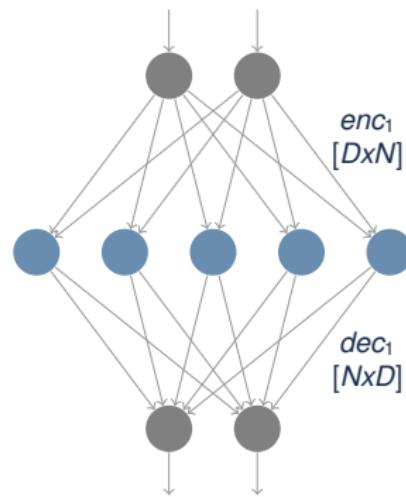
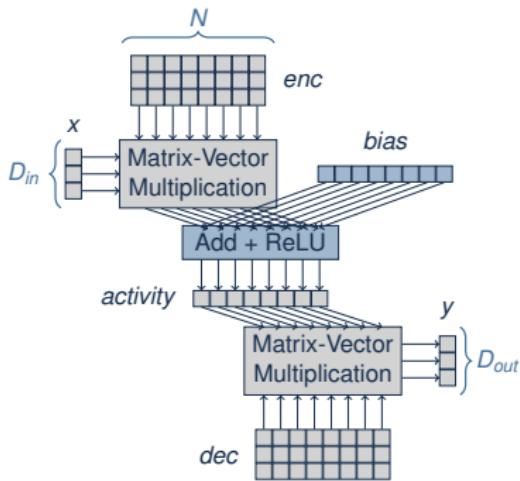
This boils down to two matrix multiplies

- ▶ One for encode



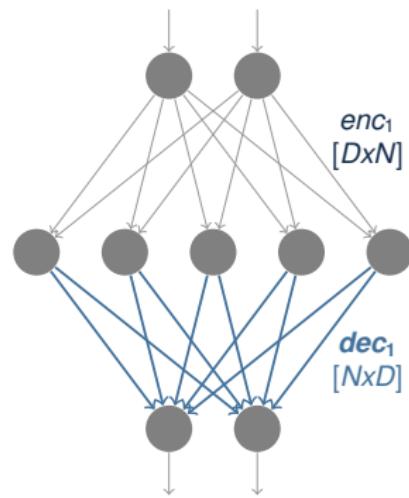
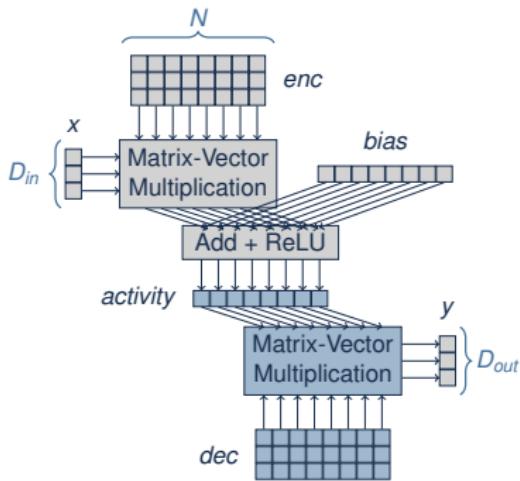
This boils down to two matrix multiplies

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



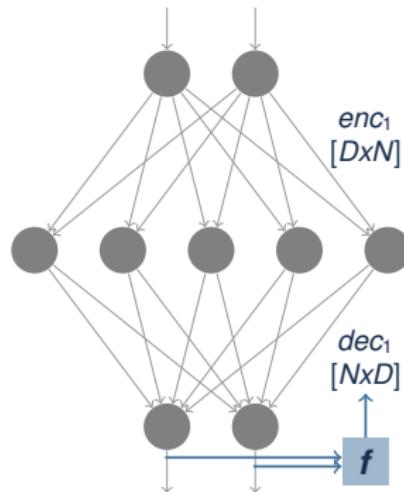
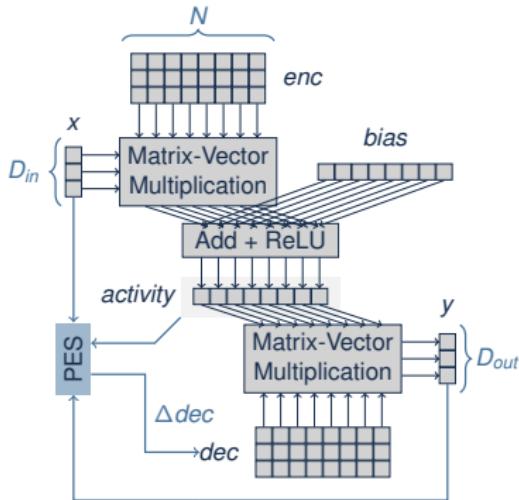
This boils down to two matrix multiplies

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



This boils down to two matrix multiplies

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



GPU is great for matrix multiplies!

- ▶ Jetson TX1 can do 1 TFLOP/s

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

In Addition,

```
1 for i < N do
2   for j < Din do
3     | ai += xj * encij
4   end
5 end
```

A single matrix multiply is simple to parallelize...

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.

The challenges we address are:

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

The challenges we address are:

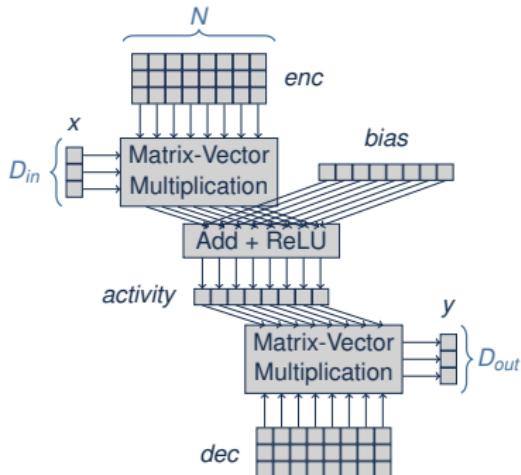
- ▶ Parallelization with HLS
  - ↳ Temporal dependency (learning) — **no batching!**
  - ↳ Structural difference between encode & decode
- ▶ Resource limitations on-chip

# 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
```

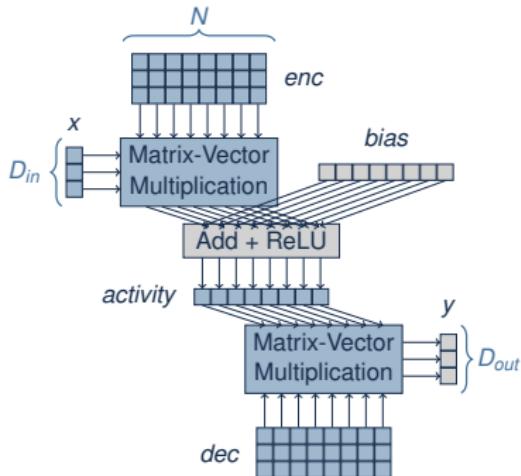


# Implementation

## 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
```



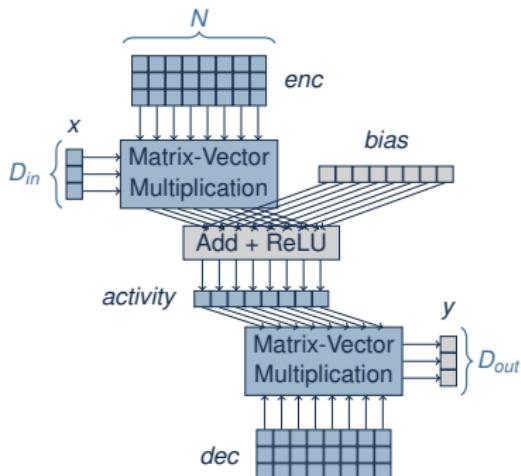
- ▶  $2 \times N \times D$

# 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
```



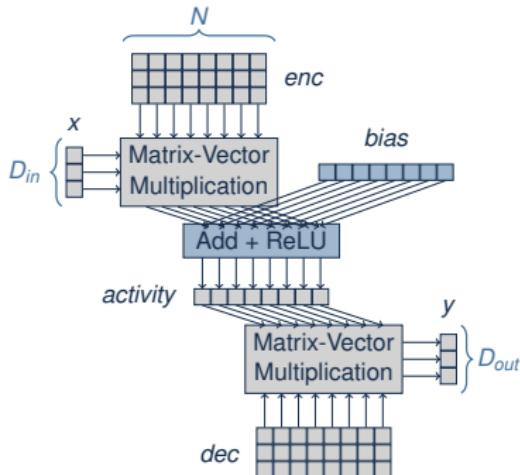
- ▶  $3 \times (2 \times N \times D)$

# 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
```



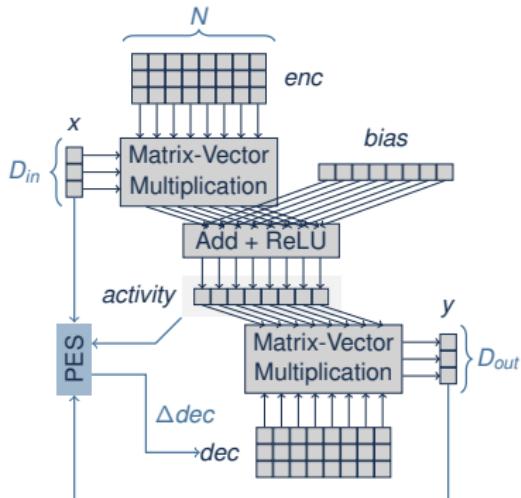
►  $3 \times (2 \times N \times D) + 2 \times N$

# 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
```



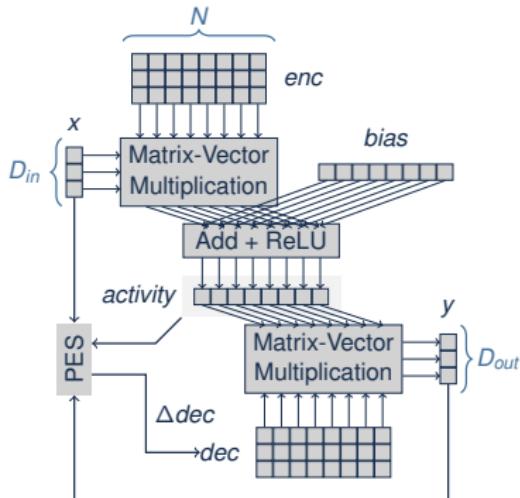
►  $3 \times (3 \times N \times D) + 2 \times N$

# 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
```



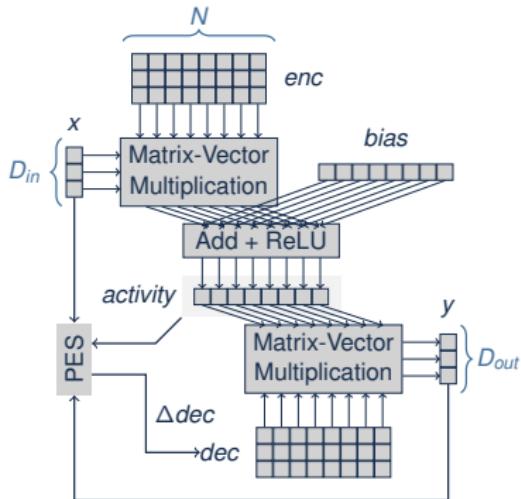
- ▶  $3 \times (3 \times N \times D) + 2 \times N$
- ▶ e.g.  $N=200$  and  $D=2$   
↳ = 4000 cycles

# 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

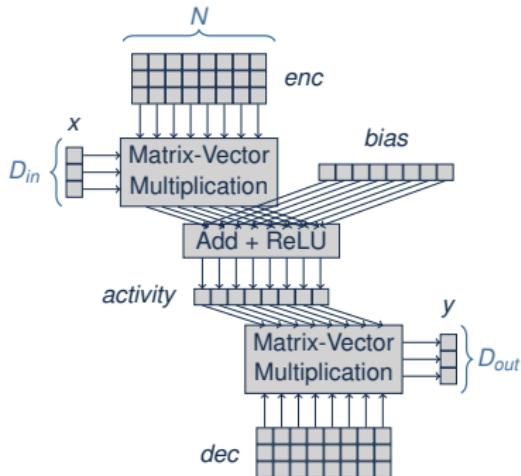
but HLS estimates **6000** cycles with naive approach!

# 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
```

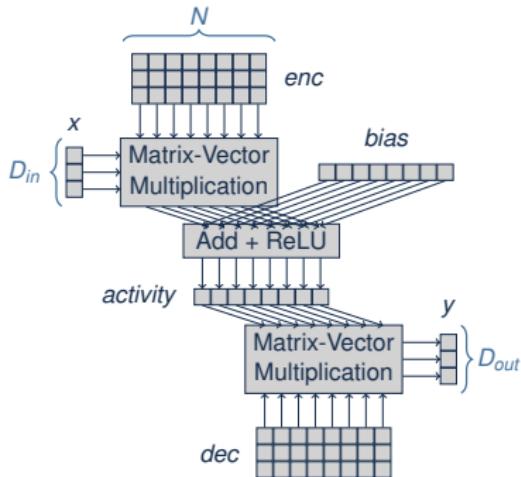


# Implementation

## 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
```

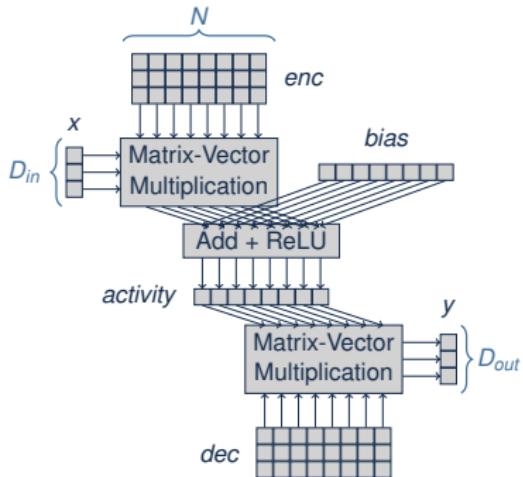


# Implementation

## 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
```



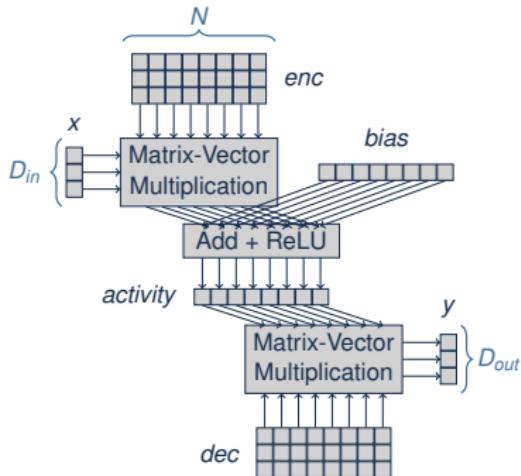
This improved II and got us to 4800 cycles.  
(recall we estimated 4000)

# 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

Now loop structure is harmonized, we want to parallelize

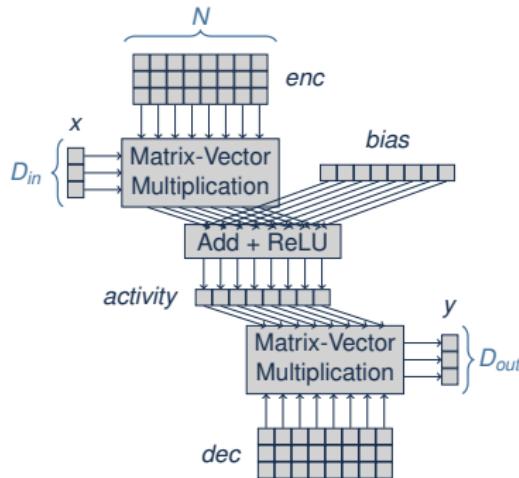
But `unroll` and `dataflow` pragmas could not automatically extract parallelism

# 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
```

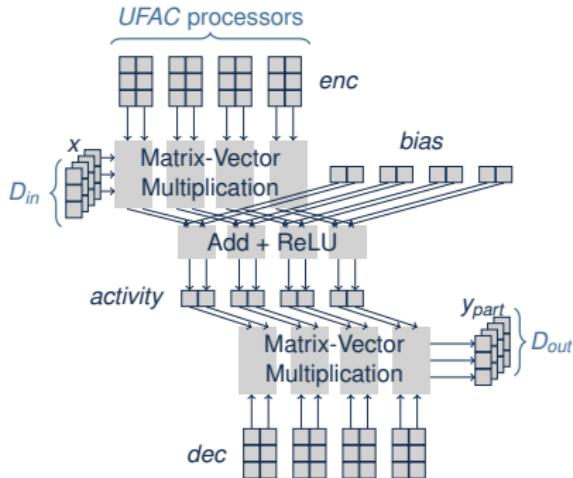


# Implementation

## High-Level Synthesis Parallelization

We add explicit parallel structure

```
1 for k < UFAC do
2   for i < ⌈ 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
```

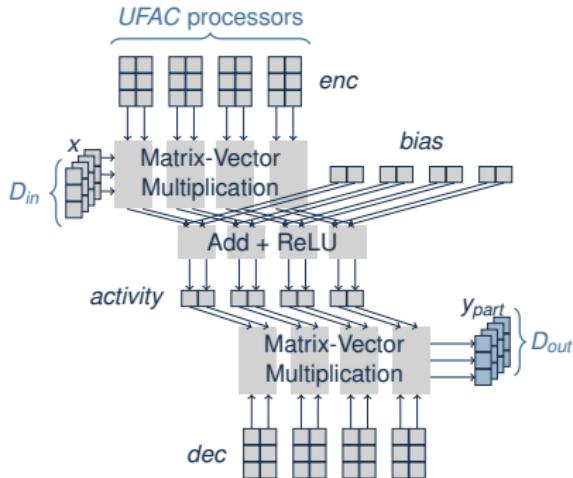


# Implementation

## 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
```

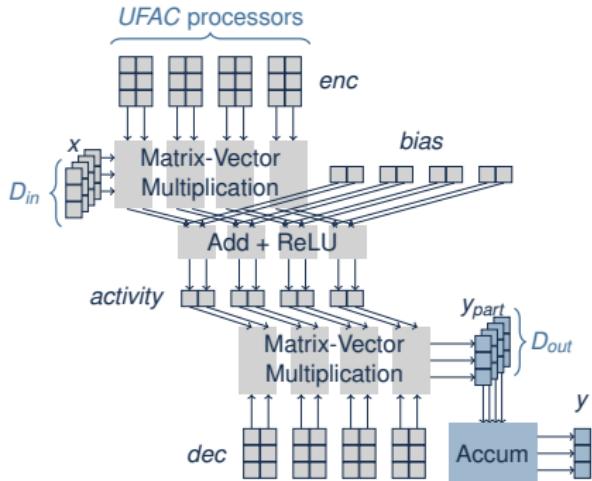


# Implementation

## 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
```

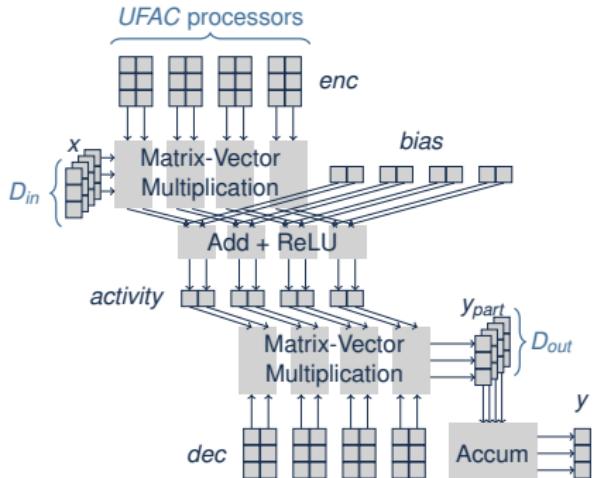


# Implementation

## 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
```

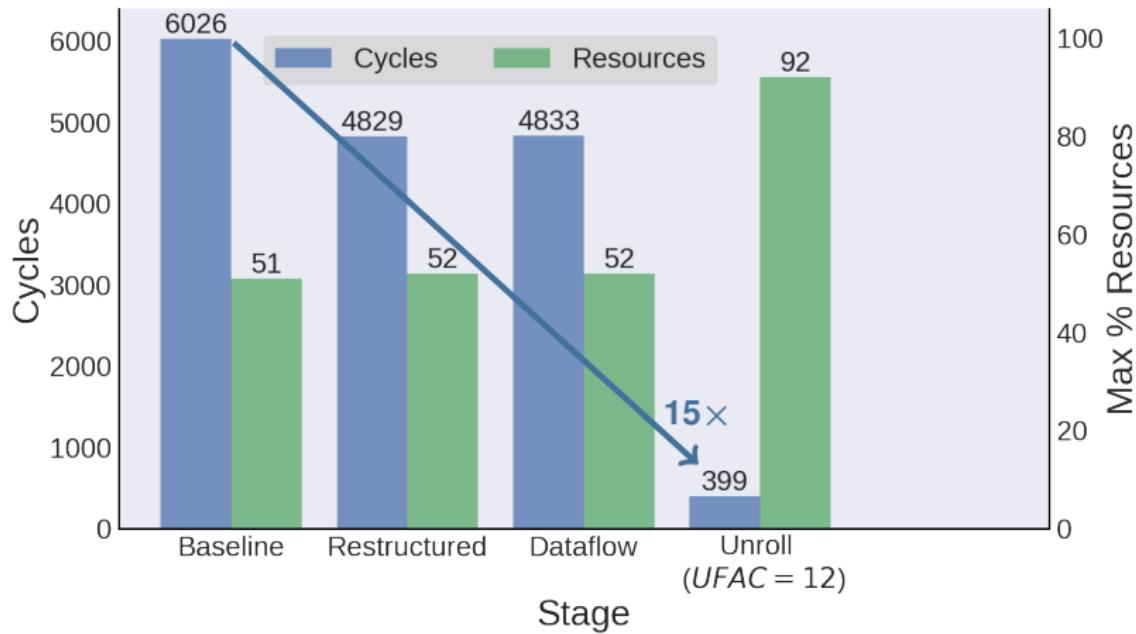


400 cycles with  $UFAC = 12$ .

# Results

## HLS Progression

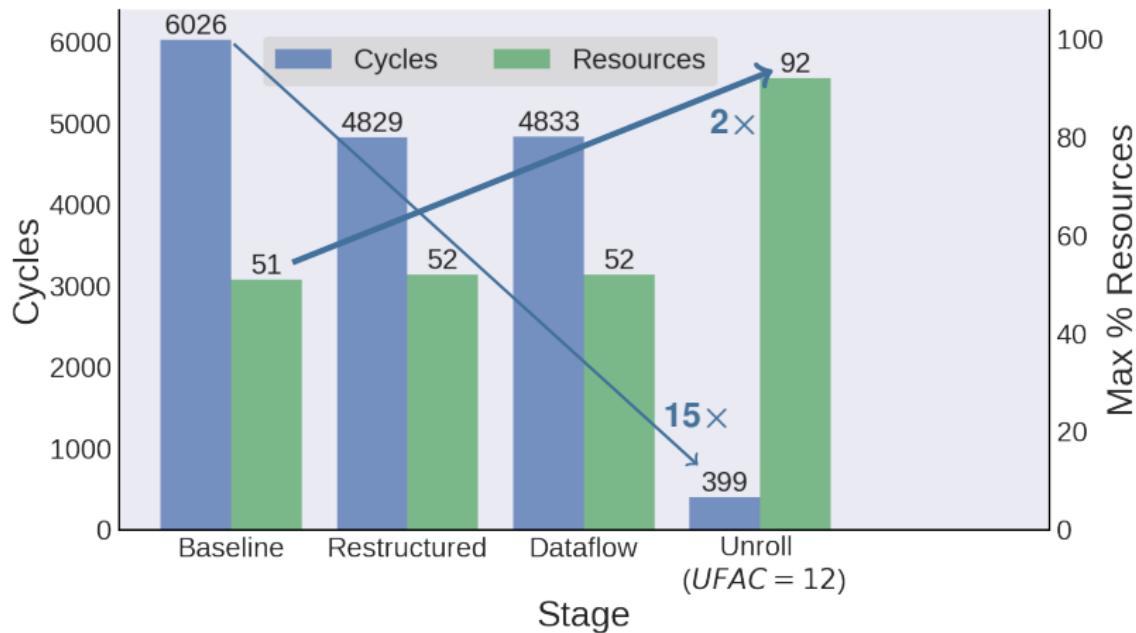
Clever HLS structure allows us to use the entire chip



# Results

## HLS Progression

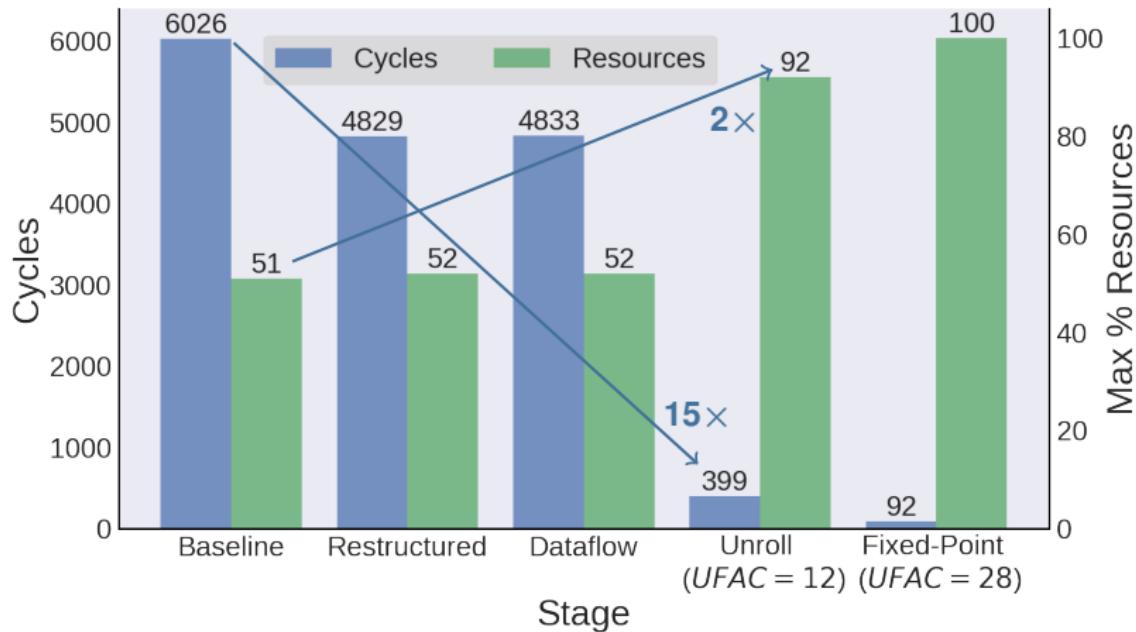
Clever HLS structure allows us to use the entire chip



# Results

## HLS Progression

### Fixed-point further shrinks compute tile



The challenges we address are:

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

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

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

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!

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

We want to minimize error and maximize performance

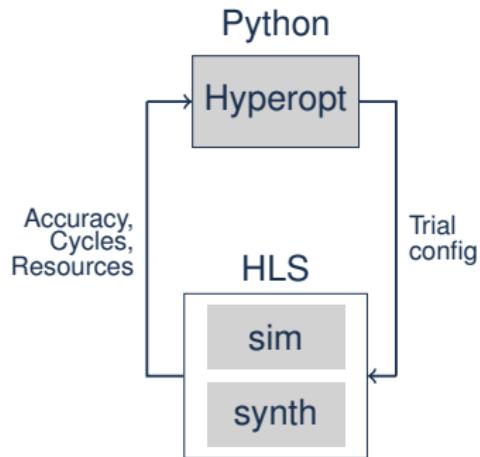
We want to minimize error and maximize performance

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

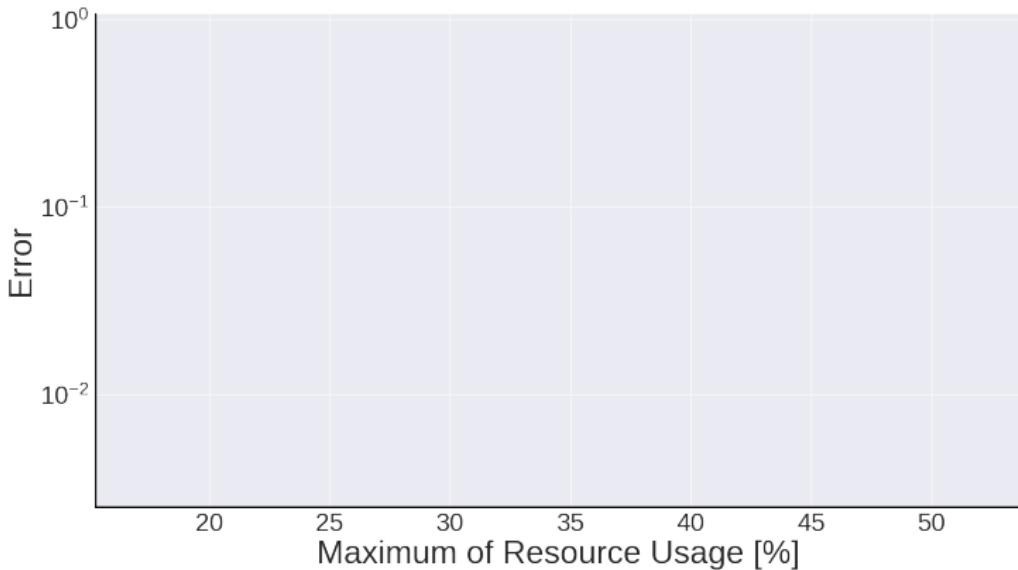
- ▶ Small, high-performance tile

Each Hyperopt trial runs:

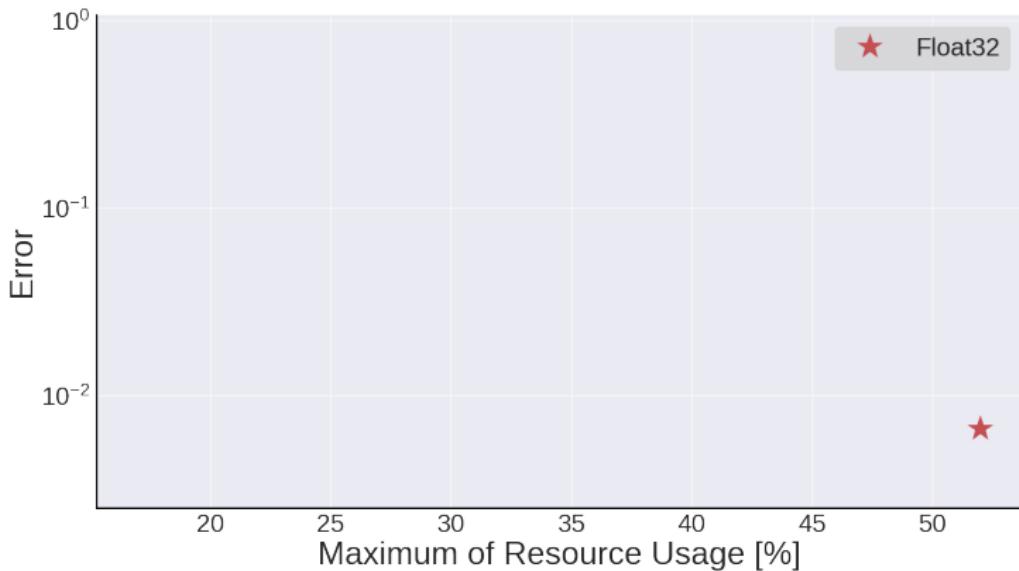
- ▶ C simulation
- ▶ C synthesis (no P & R)



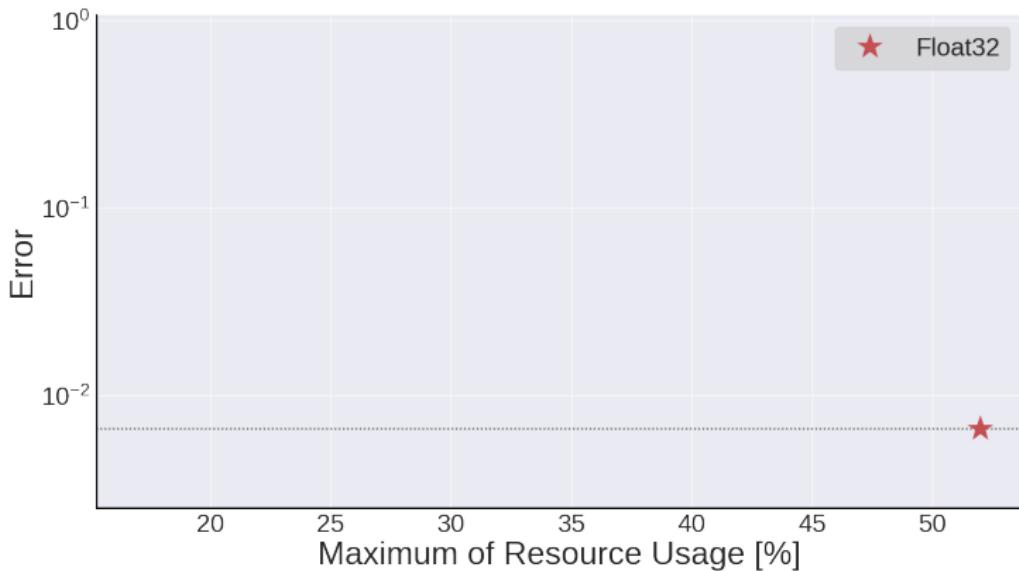
Let's look at some results:



Let's look at some results:



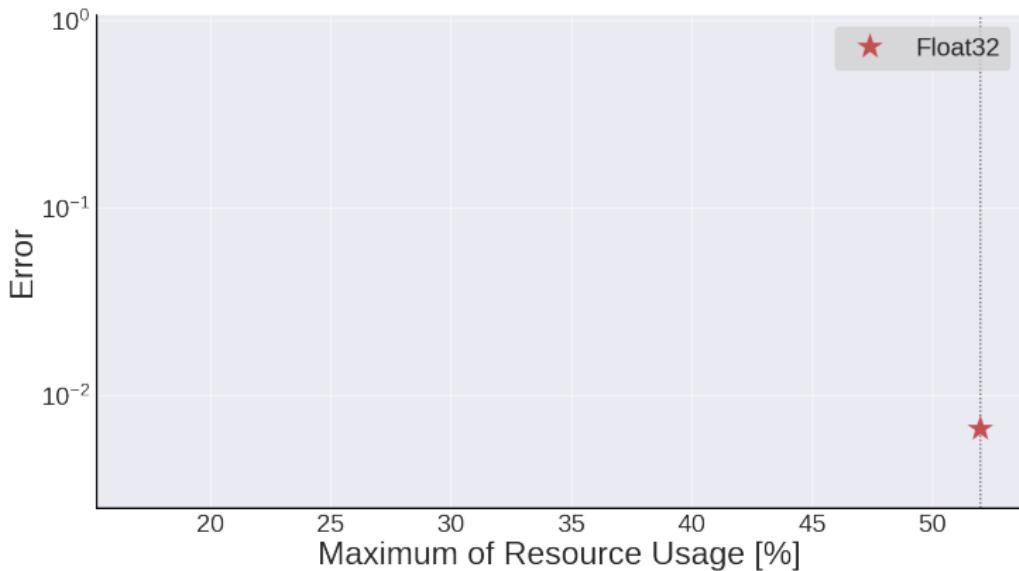
Let's look at some results:



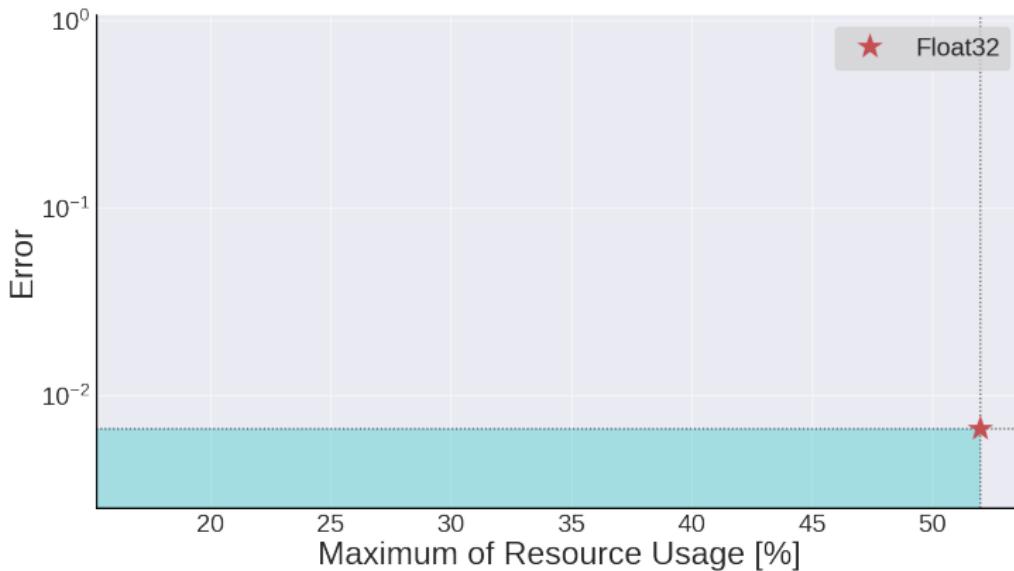
# Results

Hyperopt

Let's look at some results:



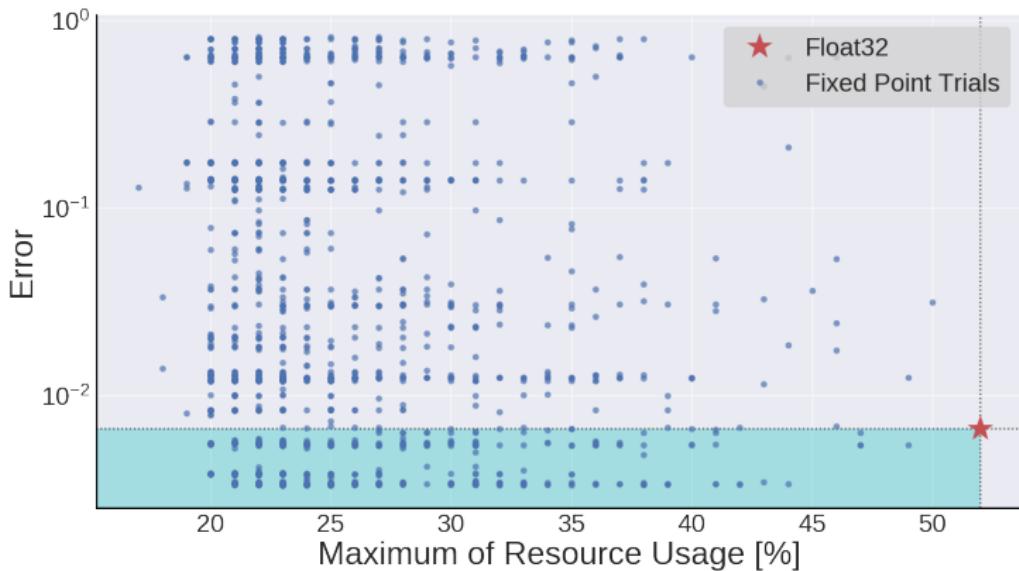
Let's look at some results:



# Results

Hyperopt

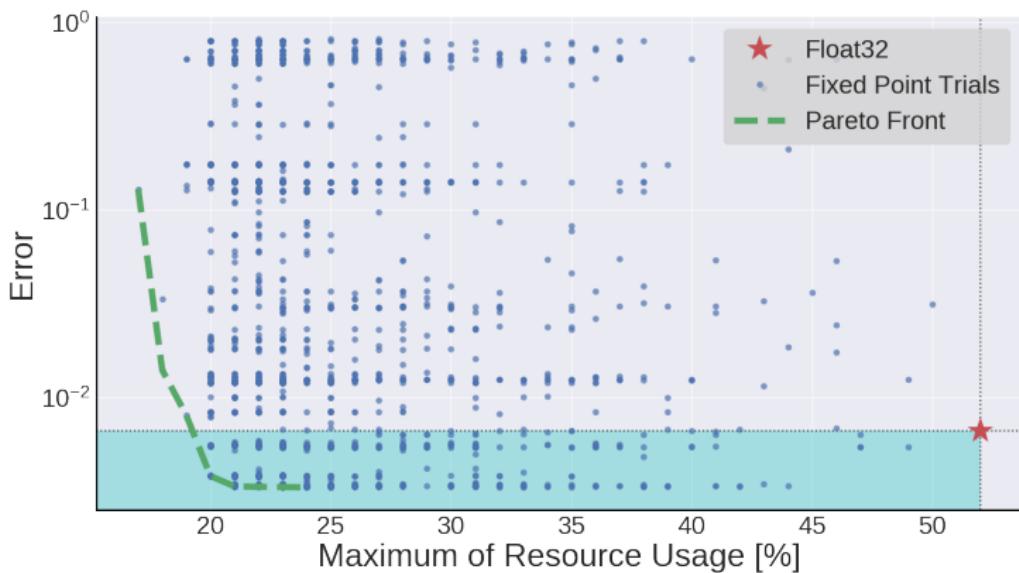
Let's look at some results:



# Results

Hyperopt

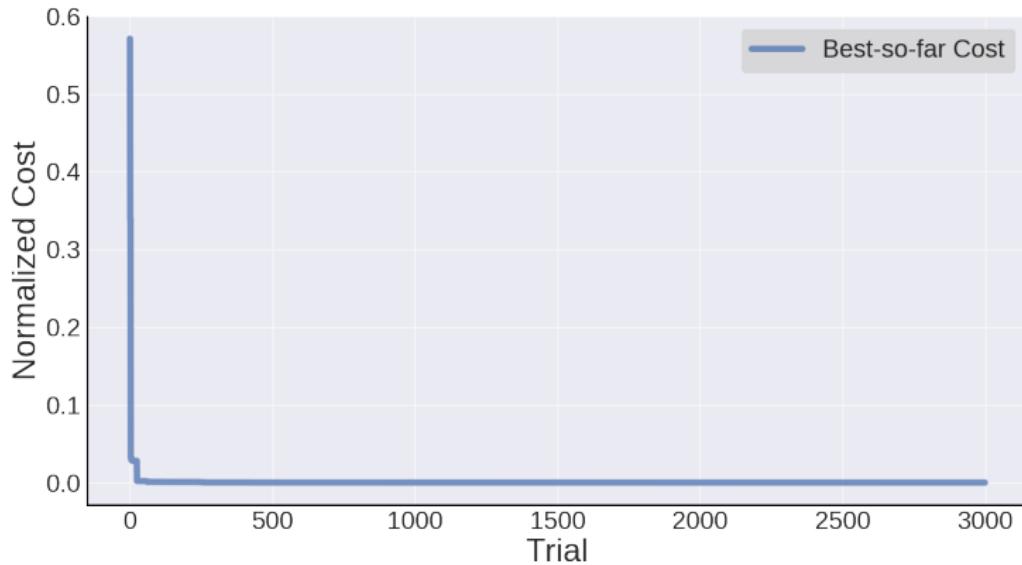
Let's look at some results:



# Results

## Hyperopt

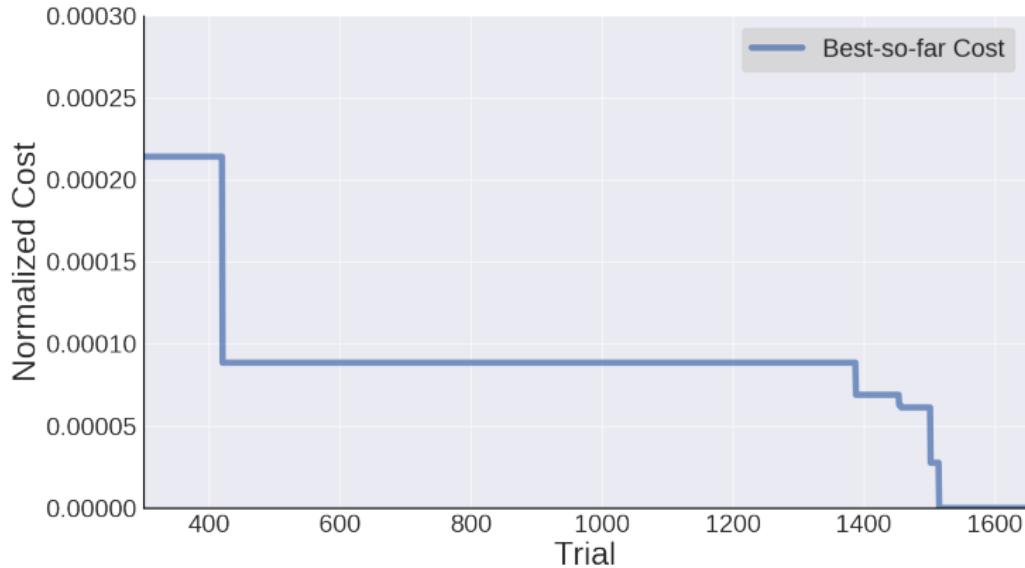
We let Hyperopt run for 3000 trials using a cost function that minimizes error and resources



# Results

## Hyperopt

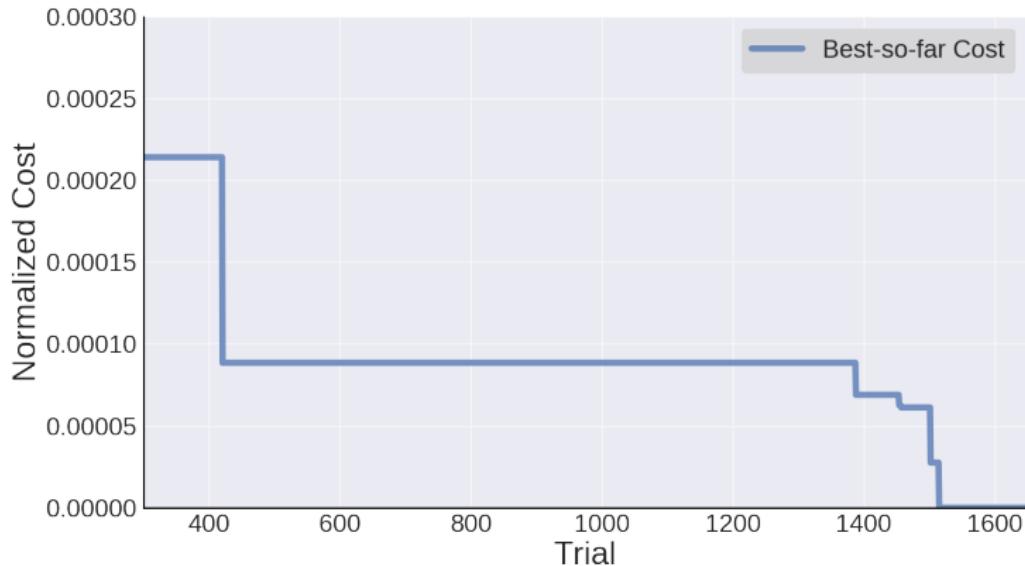
In fact, Hyperopt converges to the best solution after 1500 trials



# Results

## Hyperopt

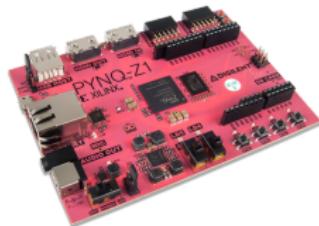
And, we are within 99% after  $\approx 400$  trials which takes 4 hours  
(i7-6700k)



# Results

Performance vs. Jetson TX1

PYNQ uses a Zynq FPGA with 28nm technology.



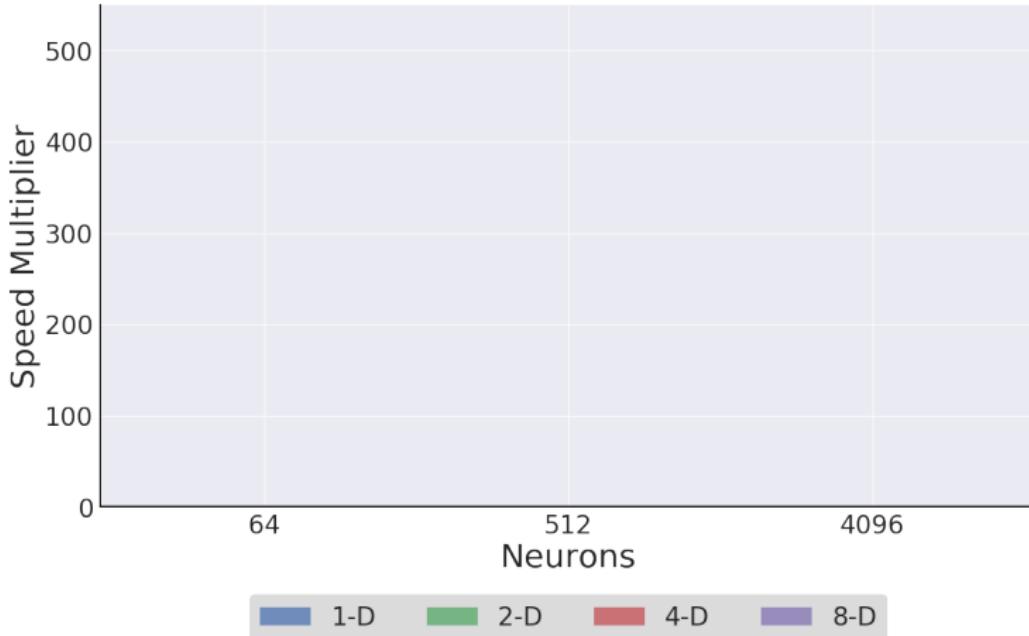
Jetson TX1 uses an embedded GPU with 20nm technology.



# Results

Performance vs. Jetson TX1

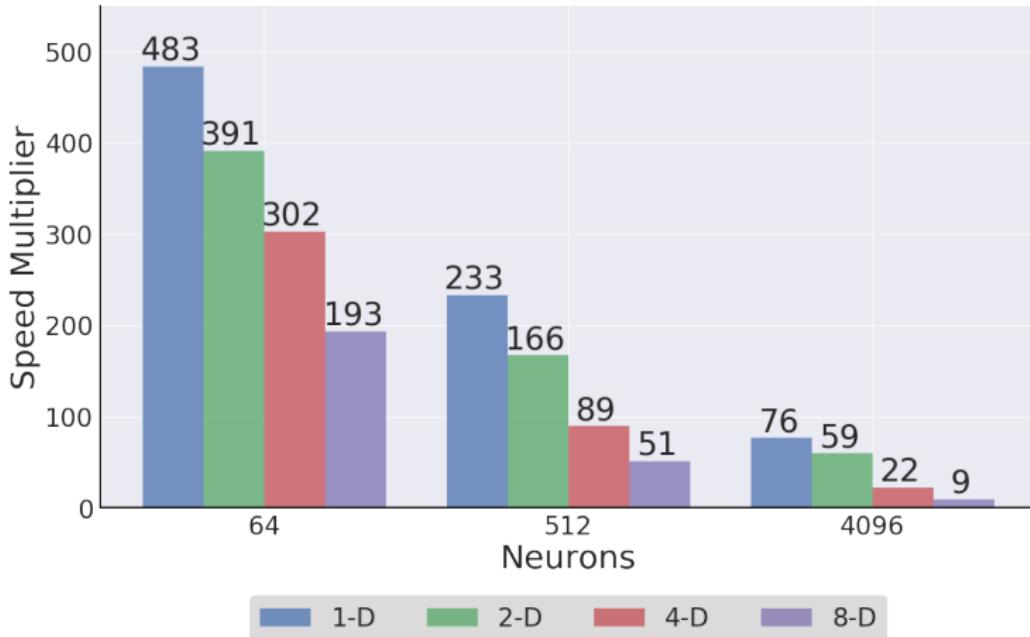
FPGA Speedup compared to cuBLAS implementation on Jetson TX1.



# Results

Performance vs. Jetson TX1

FPGA Speedup compared to cuBLAS implementation on Jetson TX1.



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]

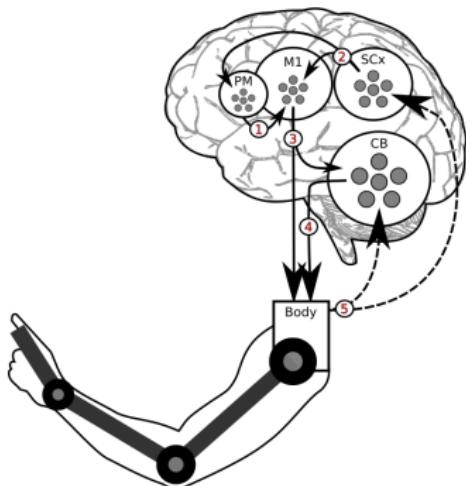
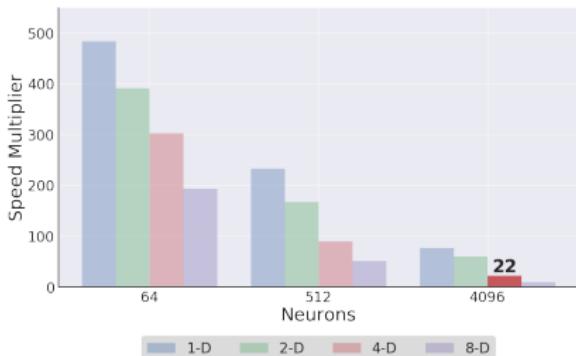


Image from "A spiking neural model of adaptive arm control" - DeWolf et al.

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

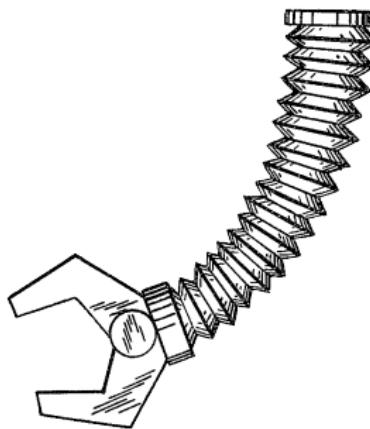
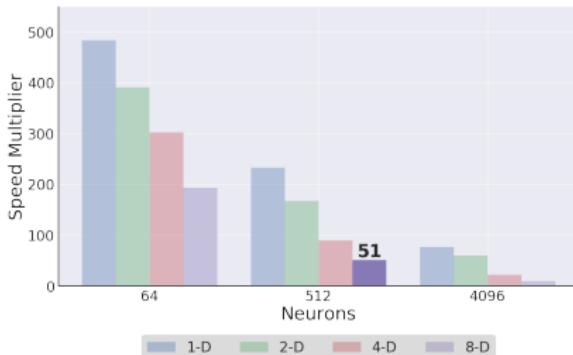


Image from Google

# Practical Applications

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

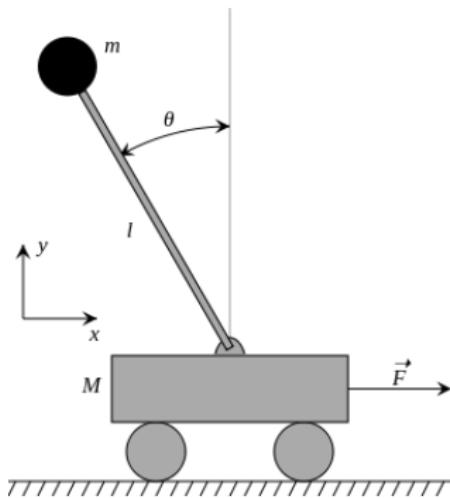
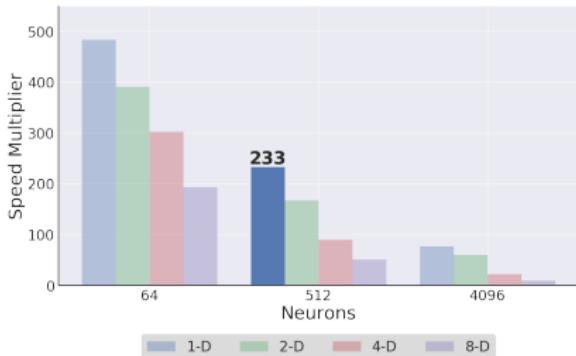


Image from Wikipedia

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

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

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
- ▶ Automatically tuned precision → 4× improvement

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 →  $4\times$  improvement
- ▶  $10\text{--}500\times$  faster than Jetson TX1

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 →  $4\times$  improvement
- ▶  $10\text{--}500\times$  faster than Jetson TX1
- ▶  $2\text{--}10\times$  less power than Jetson TX1

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

- ▶ Multiple neuron types
- ▶ Larger networks
- ▶ Flexible connectivity
- ▶ Better CNN support

**NengoFPGA** is available for purchase from  
**Applied Brain Research** (url: [abr.rocks](http://abr.rocks))