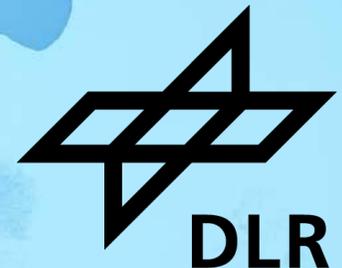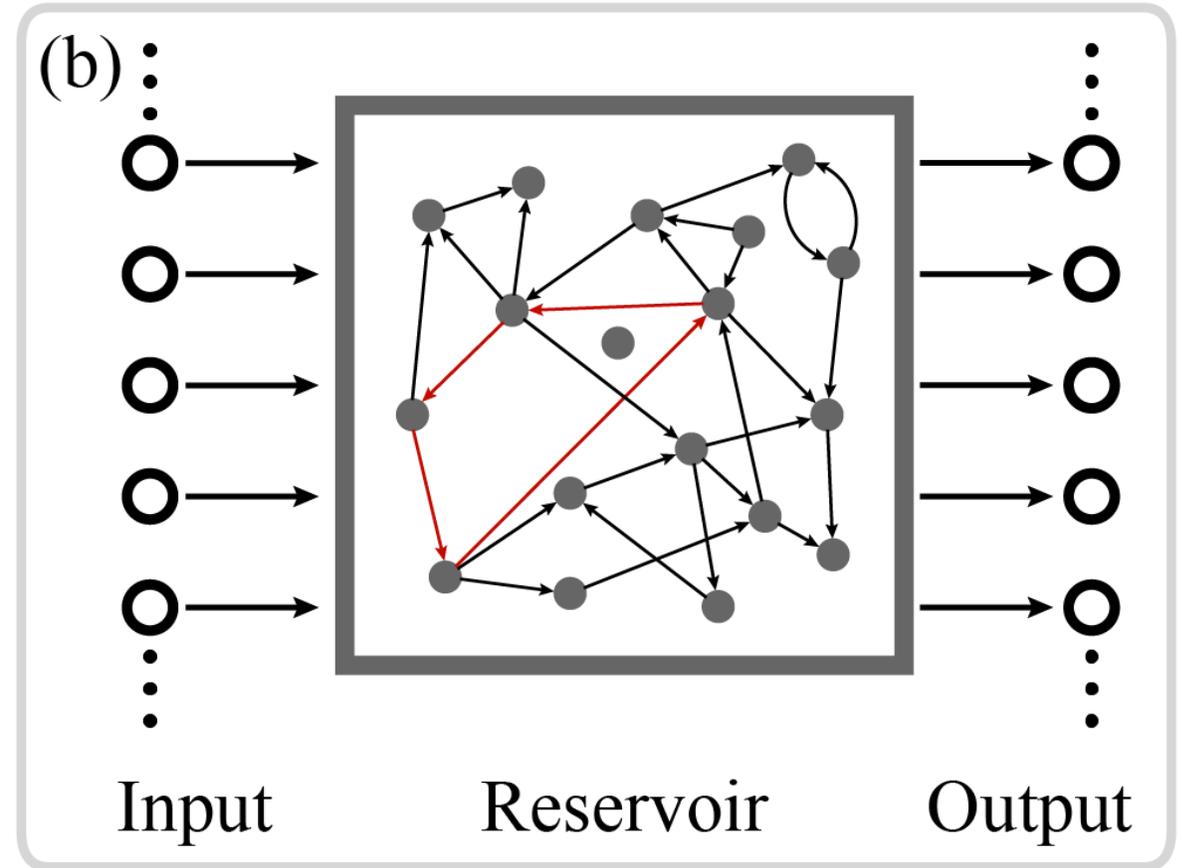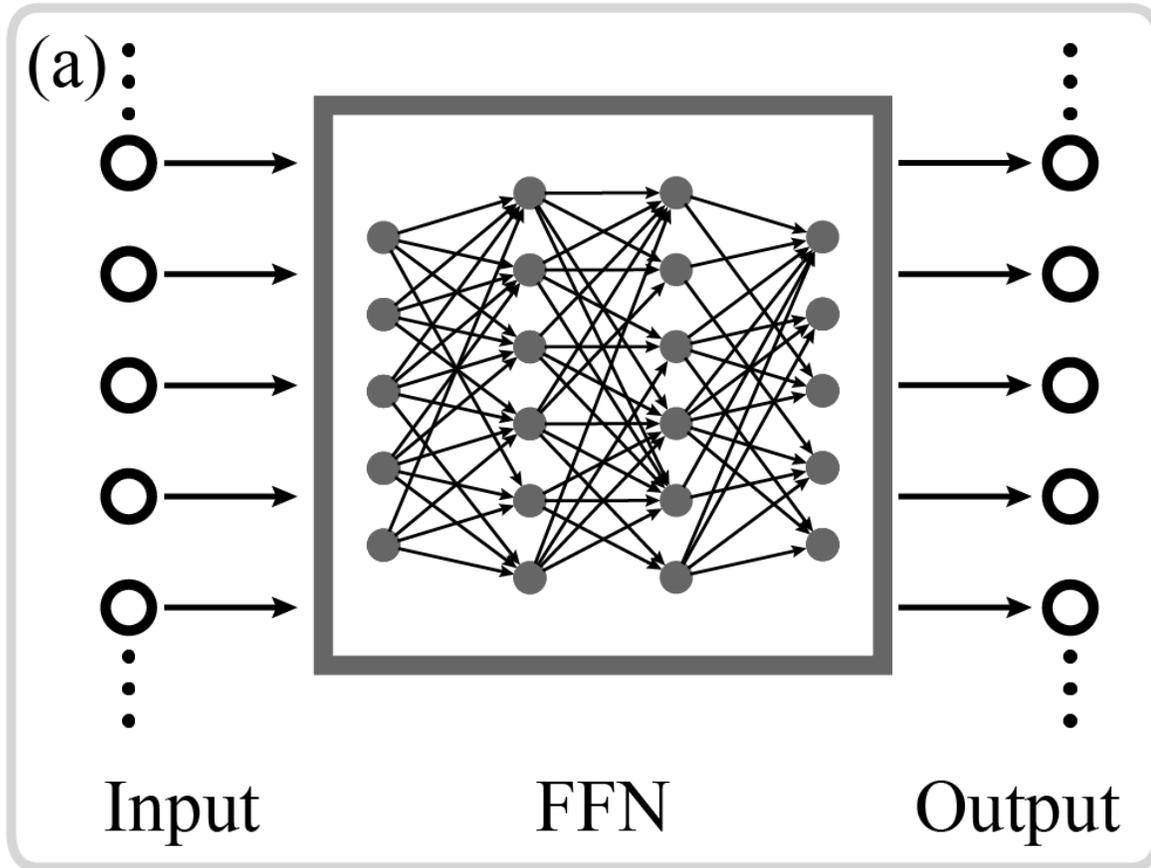# FPGA IMPLEMENTATION OF NGRC FOR PREDICTING DYNAMICAL SYSTEMS

**By João Folhadela and Daniel Köglmayr – DLR KI-EXE**

**2nd FPGA Developers' Forum (FDF) meeting (CERN)**

João C.L. Folhadela, Institute for AI Safety and Security, 21.05.2025

DLR

João C.L. Folhadela, Institute for AI Safety and Security, 21.05.2025

Optimize $\boldsymbol{W}_{out}$ via Ridge regression. Minimize loss:

$$\sum_{t=-T}^{0} \| \boldsymbol{W}_{out}\boldsymbol{r}(t) - \boldsymbol{y_R}(t) \|^2 + \beta \| \boldsymbol{W}_{out} \|^2$$

$\boldsymbol{y_R}(t)$: known real output

$\beta$: regularization parameter

$\boldsymbol{W}_{out}$ is then obtained via:

$$\boldsymbol{W}_{out} = (\boldsymbol{r}^{\intercal}\boldsymbol{r} + \beta\,\boldsymbol{I})^{-1}\,\boldsymbol{r}^{\intercal}y$$

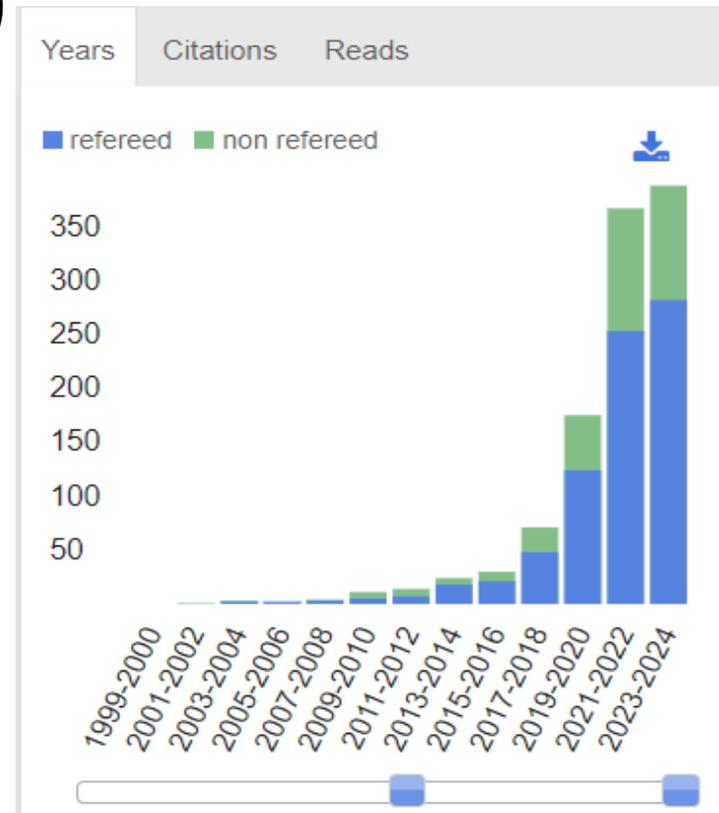An analytic solution is usually much faster to compute than gradient descent.
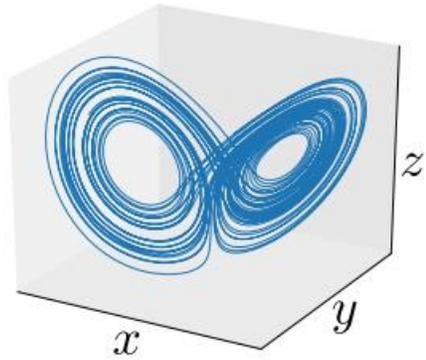
# RC SOTA
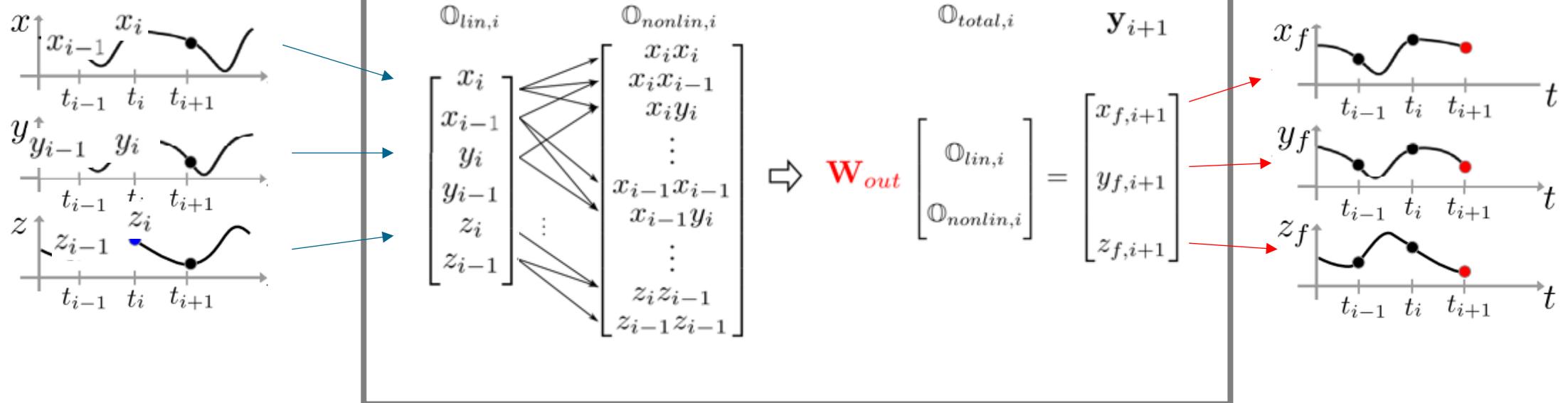


All methods have two basic architectural components:
- Nonlinear dimensionality expansion of the input data
- Regression based mapping on the target data

João C.L. Folhadela, Institute for AI Safety and Security, 21.05.2025

# Introduction – Next Generation Reservoir Computing (NGRC)

- Efficient
- Less training data required

João C.L. Folhadela, Institute for AI Safety and Security, 21.05.2025

D. Gauthier et al., Nat. Comm., 2021

# Introduction – NGRC on Lorenz Attractor

- ▪ We train with Python's numpy's linalg solver   $W_{out} = (r^{\mathsf{T}}r + \beta\, I)^{-1}\, r^{\mathsf{T}} Y_{train}$

- ▪ Some parameters:

| Train Steps | k | β | $W_{out}$ shape |
|:---:|:---:|:---:|:---:|
| 400 | 2 | $9 \times 10^{-2}$ | 27x3 |

k – timesteps β - regularization parameter

- ▪ Then inference is performed on FPGA
  - ▪ Transfer $W_{out}$ and first linear state to FPGA memory
  - ▪ Prediction loop

# Implementation – The NGRC prediction loop

1.  Creating the Reservoir
    - Read linear state from memory – current and previous timestep
    - Non-linear expansion – <span style="color:red">multiplying and powering the linear state</span>

2.  Predicting
    - Matrix multiplication, several vector dot products – <span style="color:red">Multiply and Accumulate</span> (MACs)
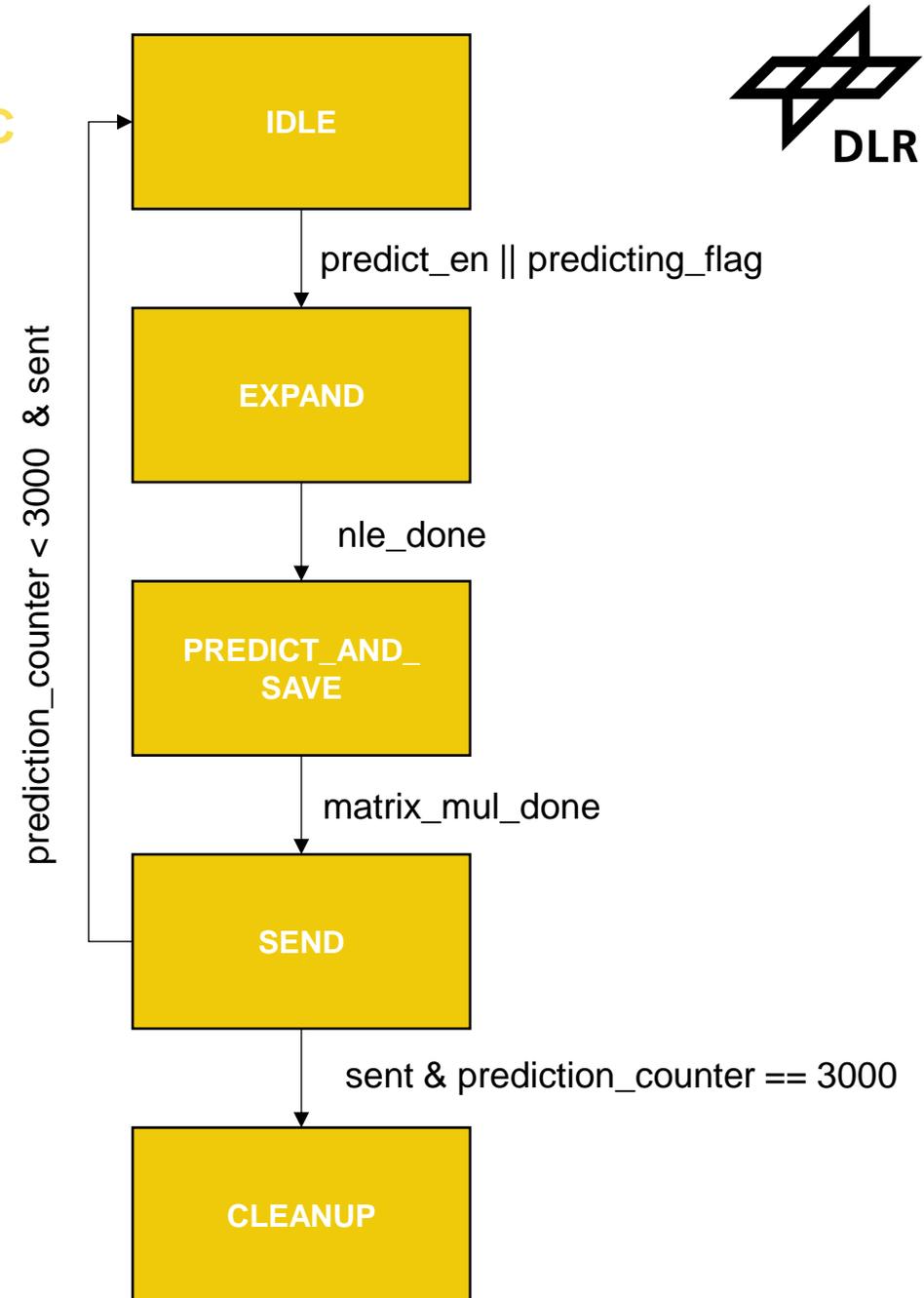
3.  Save prediction


<span style="color:red">In red</span> – arithmetics are performed in float 32 with Vivado floating-point IP

# Implementation – The FSM

- ## IDLE
  - Sets linear state
  - Waits for prediction loop to begin/continue
- ## EXPAND
  - Performs the non-linear expansion
- ## PREDICT_AND_SAVE
  - Calculates predicton
  - Saves it
- ## SEND
  - Send the current prediction via UART
- ## CLEANUP
  - Resets all counters and variables

NGRC

João C.L. Folhadela, Institute for AI Safety and Security, 21.05.2025
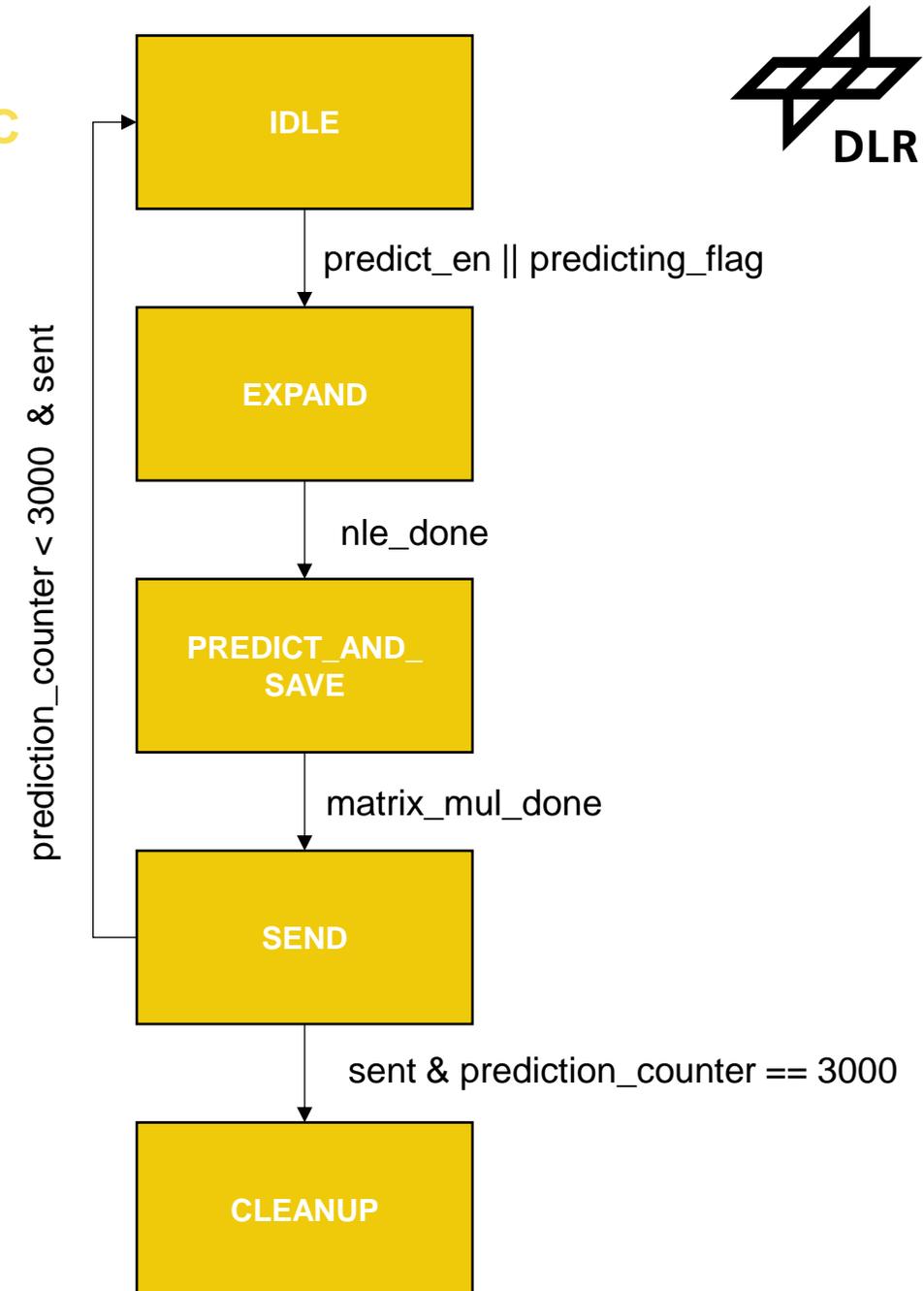
# Implementation – The FSM

EXPAND

- Implemented as a module that expands a vector of size 6 into size 27 (6+21)

- Each multiplication is computed in parallel, with 32-bit float numbers:
  - Vivado Floating-point IP
    - *Multiply* operation
    - *Full DSP usage* (21 multiplications x 2 DSP = 42 DSPs)
    - Maximum Latency used (8 clock cycles)

- Concatenate the input vector with each of the multiplication results.

# Implementation – The FSM

- IDLE
  - Sets linear state
  - Waits for prediction loop to begin/continue
- EXPAND
  - Performs the non-linear expansion
- PREDICT_AND_SAVE
  - Calculates predicton
  - Saves it
- SEND
  - Send the current prediction via UART
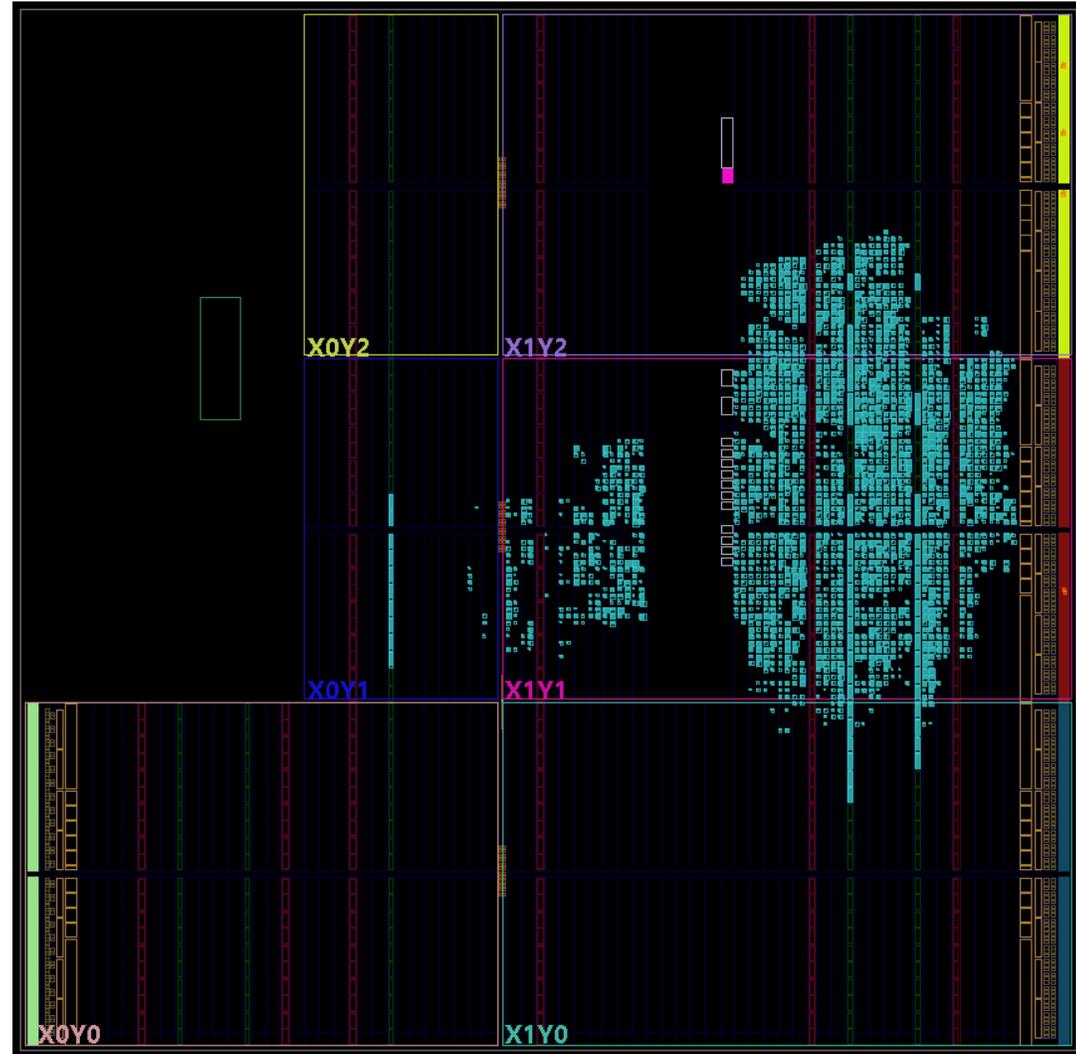- CLEANUP
  - Resets all counters and variables

NGRC

```
        ┌──────────────┐
        │     IDLE     │
        └──────┬───────┘
               │ predict_en || predicting_flag
        ┌──────▼───────┐
        │    EXPAND    │
        └──────┬───────┘
               │ nle_done
   ┌──────────▼─────────┐
   │  PREDICT_AND_SAVE  │
   └──────────┬─────────┘
              │ matrix_mul_done
        ┌─────▼────────┐
        │     SEND     │
        └─────┬────────┘
              │ sent & prediction_counter == 3000
        ┌─────▼────────┐
        │   CLEANUP    │
        └──────────────┘
```

prediction_counter < 3000 & sent

# Implementation – The FSM

PREDICT_AND_
SAVE

DLR

- Perform matrix multiplication of the reservoir state (a vector) with the weights, $W_{out}$:
  - 3 MACs implemented in parallel with Vivado Floating-point IP that compute:
    - *Fused Multiply-Add* operation
    - *Full DSP usage* (3 MACs x 4 DSP = 12 DSPs)
    - 10 clock cycles latency used
  - 10 latency x 27 accumulations = 270 clock cycles in total
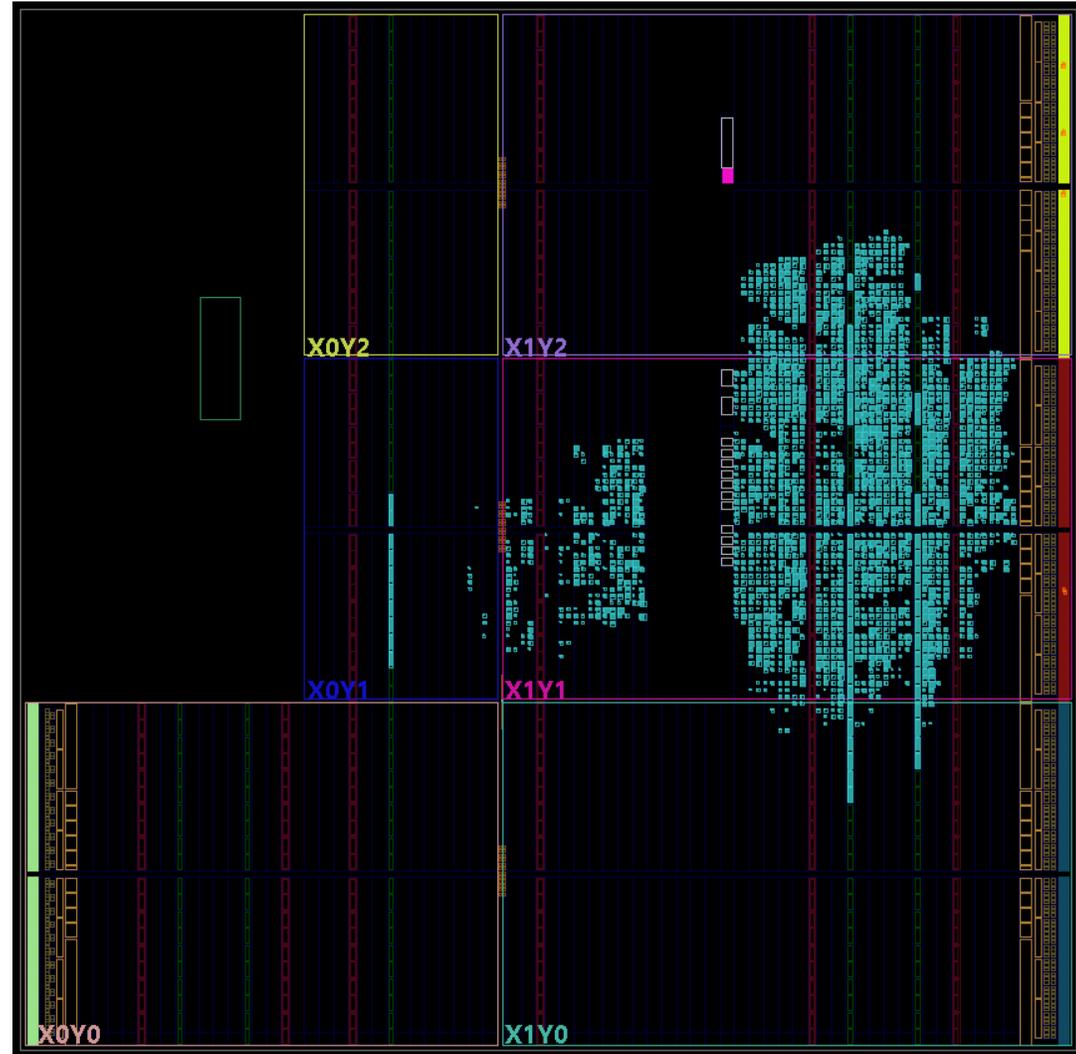- Save prediction (to update next linear state)

# Resource usage

| Resource | Used / Available |
|----------|------------------|
| LUT | 5,132 / 53,200 |
| LUTRAM | 270 / 17,400 |
| FF | 6,475 / 106,400 |
| DSP | 54 / 220 |

# Other reports

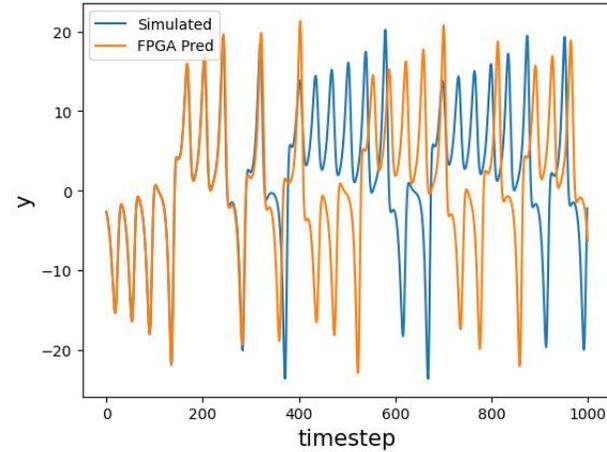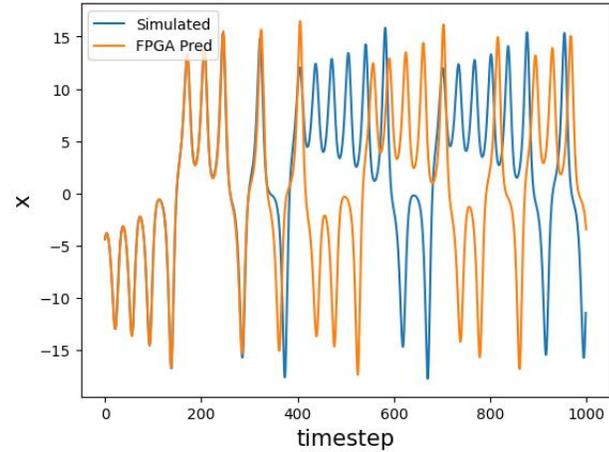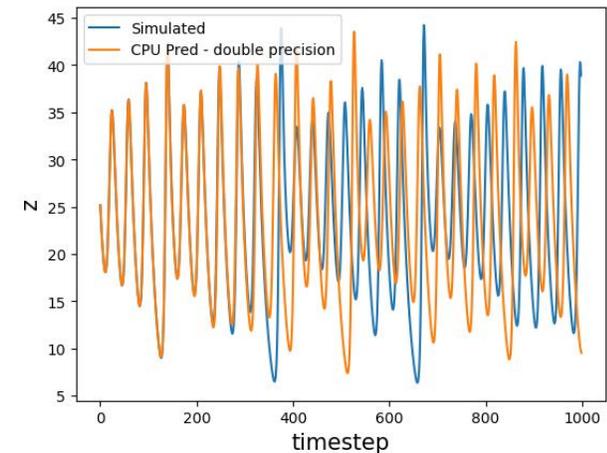| Property | Value |
|---|---|
| On-Chip Power | 0.294 W |
| Dynamic | 0.186 W |
| Static | 0.109 W |
| Clock Speed | 125 MHz |
| Worst Negative Slack | 1.3 ns |

# Evaluation – Simulated vs. Predicted Lorenz

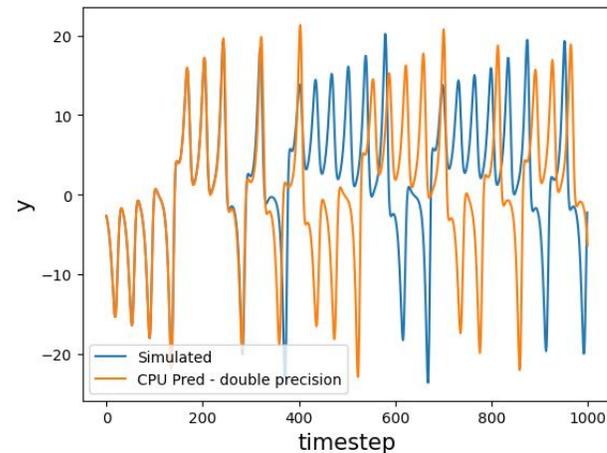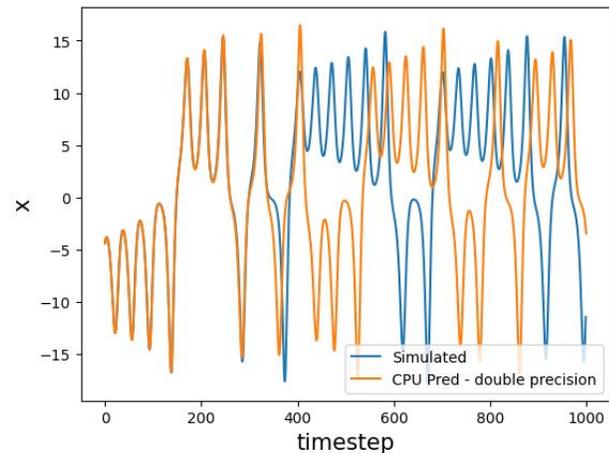Prediction over 3000 timesteps

# Evaluation – Simulated vs. Predicted Lorenz

## On FPGA



## On CPU – with double precision

João C.L. Folhadela, Institute for AI Safety and Security, 21.05.2025

# Latency Evaluation – Simulated vs. Predicted Lorenz



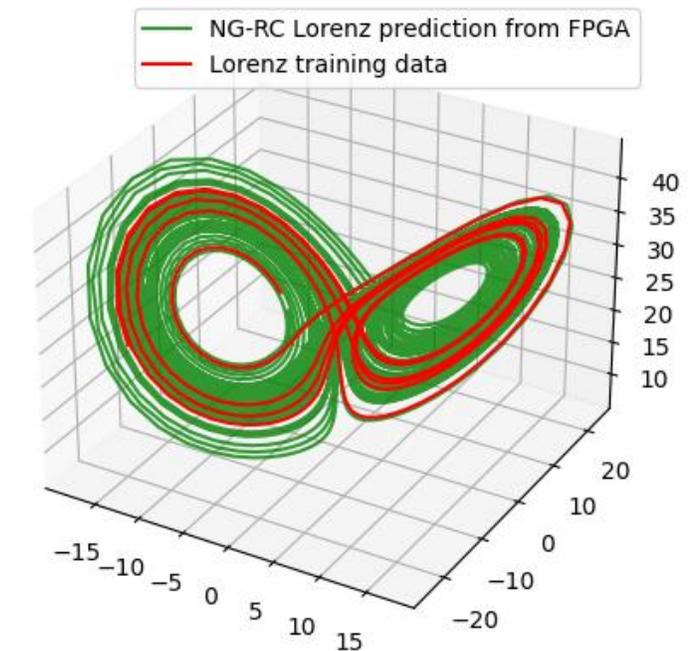| | FPGA | CPU |
|---|---|---|
| Device | Zynq7000-20 | i7-1365U 13[th] Gen |
| Clock Speed | 125 MHz | 4.8 GHz |
| Inference time | 2.5 µs | 50 µs |

- Easy to implement

- Low power
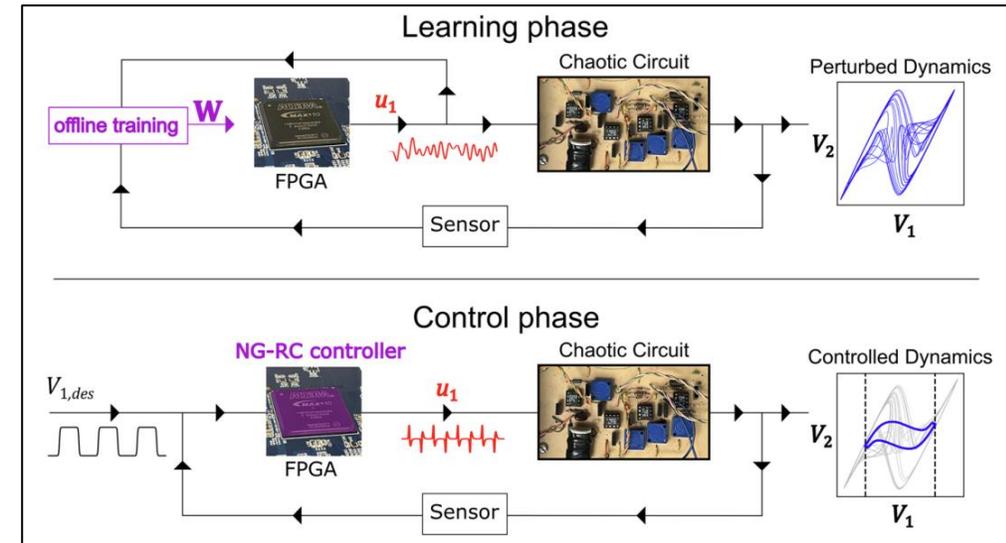
- Fast

**And real applications?**

# Applications of NGRC



- **Control of chaotic circuit dynamics**

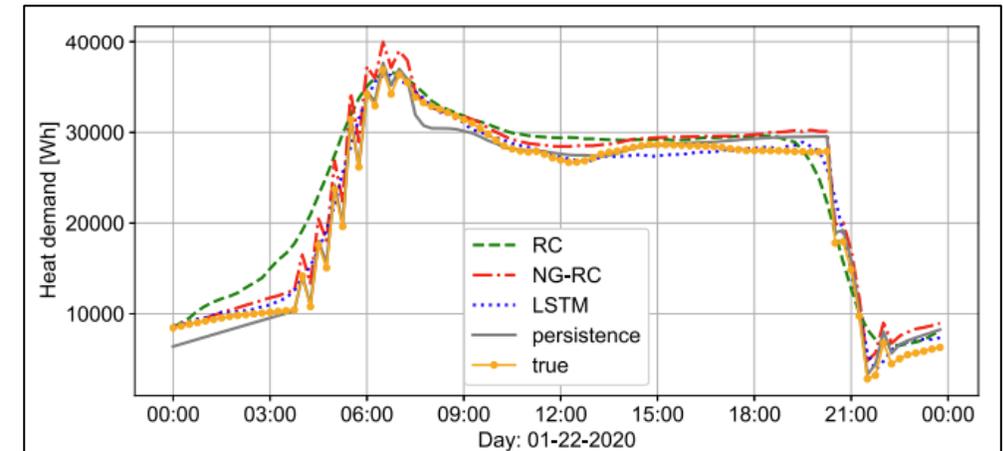Kent, R.M. et. al *Controlling chaos using edge computing hardware*

Nat Commun 15, 3886 (2024)

- **Energy demand forecasting**

Karoline Brucke et. al *Benchmarking reservoir computing for residential energy demand forecasting*

Energy and Buildings 314 (2024)



- **Financial markets, weather, costal line erosion, robotic control, etc.**

# THANK YOU!

Questions?

João C.L. Folhadela, Institute for AI Safety and Security, 21.05.2025

DLR