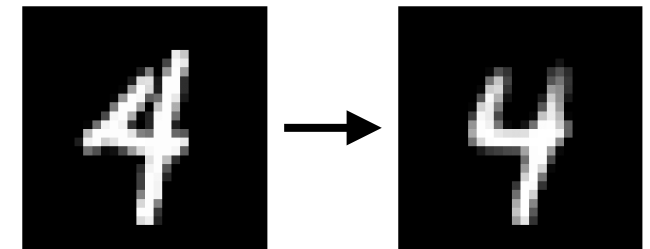
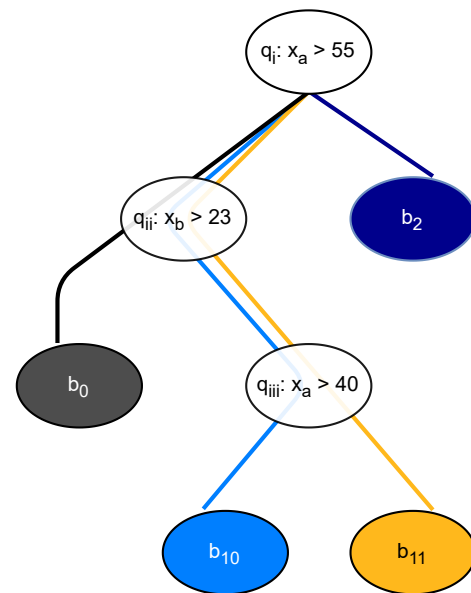


Tree-distilled

# Autoencoder on FPGA to anomaly detect compress data

^



Tae Min Hong



University of  
Pittsburgh

AI4EIC, MIT

October 27, 2025

<https://indico.bnl.gov/event/28082/contributions/115514/>



# Outline

## Introduction

## Anomaly Detection

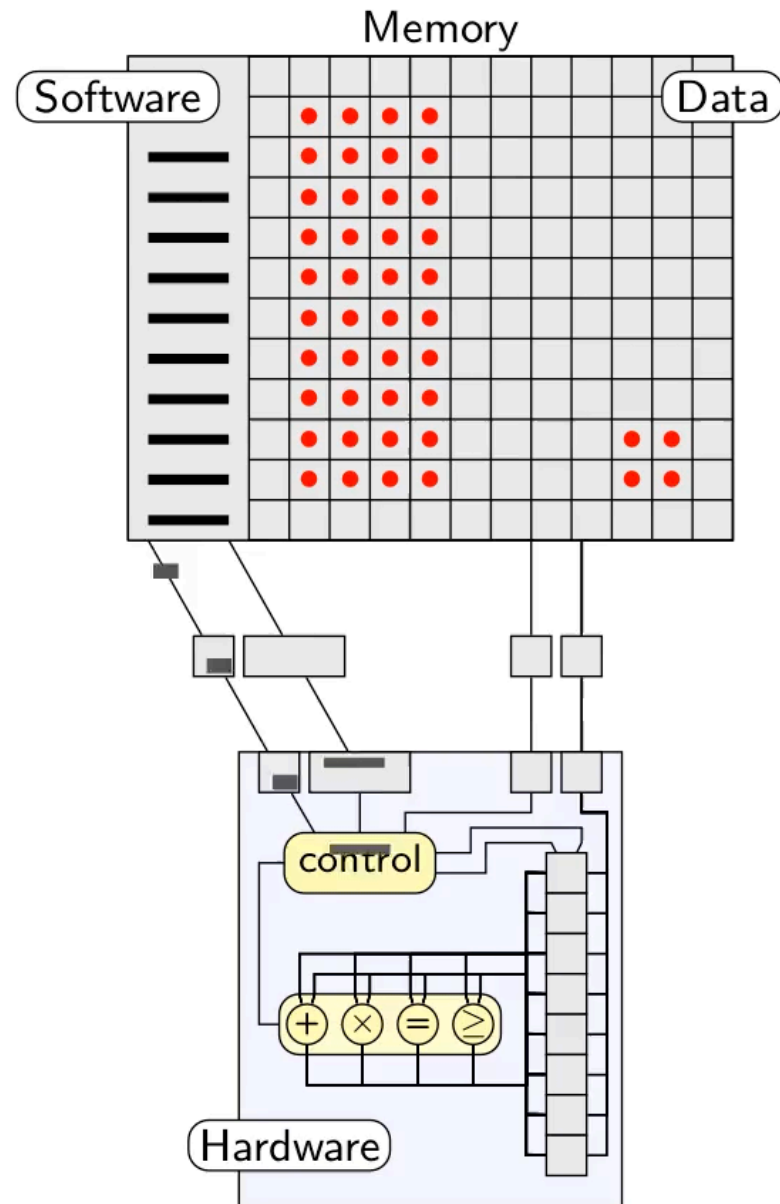
## Trees on FPGA

## Data Compression

- **FPGA** Field programmable gate arrays
- **AE** Autoencoder
- **DT-AE** Decision Tree AE
- Train DT-AE → Deploy Trees
- Regression
- Decision Tree structure
- HLS vs. VHDL  1st time showing
- Train VAE → Deploy Trees
- Data compression  1st time showing
- Code & git & slides & videos

# Why FPGA (vs. CPU)

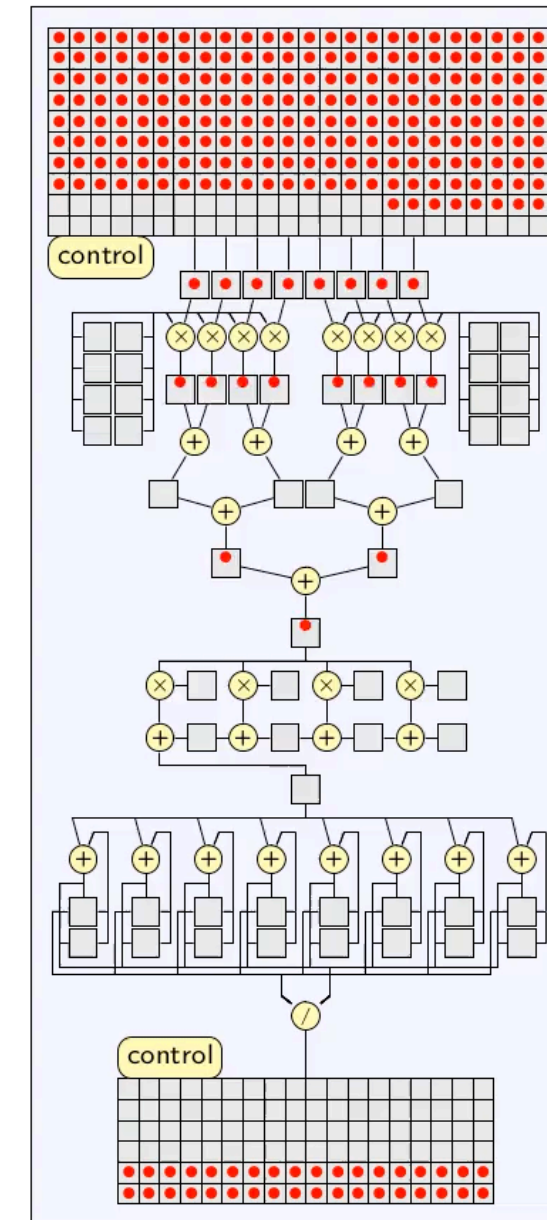
Animation from [https://qbaylogic.com/wp-content/uploads/2024/05/Animation\\_CPU-vs-FPGA.mp4](https://qbaylogic.com/wp-content/uploads/2024/05/Animation_CPU-vs-FPGA.mp4)



**Normal processor**

Software controls computation process

One operation per clock tick  
typical timing  $10^{-3}$  ms



**FPGA**

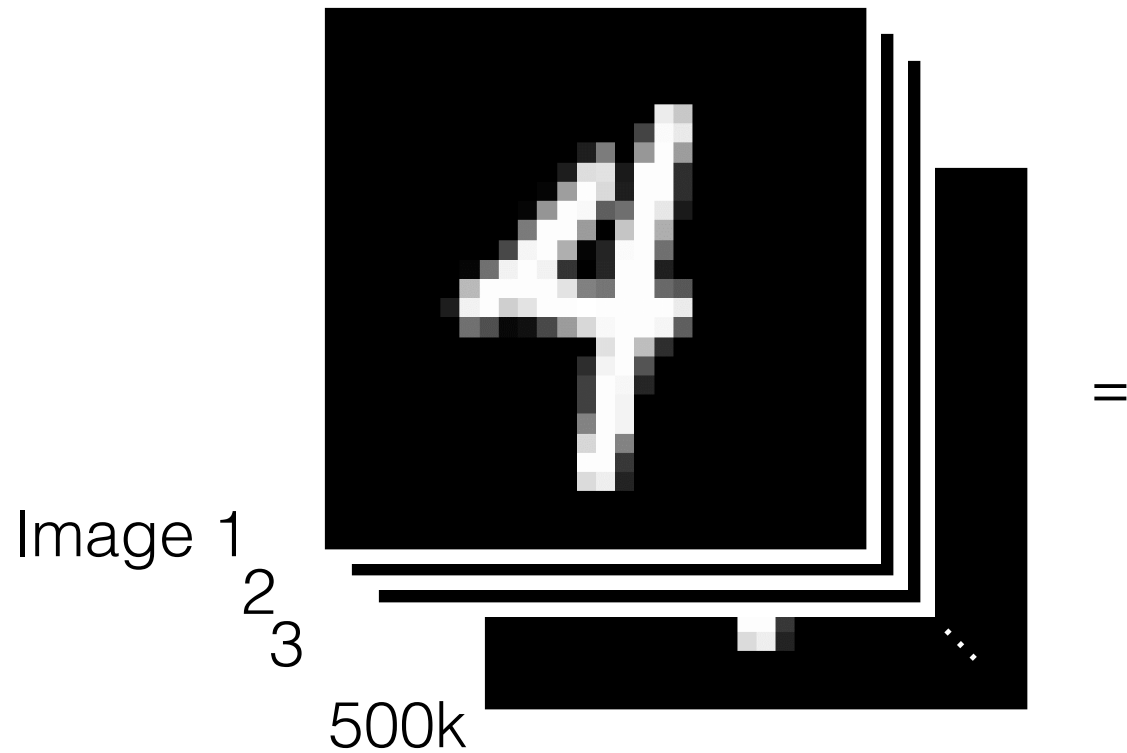
Computation integrated in hardware

Many operations per clock tick  
typical latency  $\ll \mu\text{s}$ , throughput  $\gg 1$  MHz

# Autoencoder recap

Example: handwritten digits

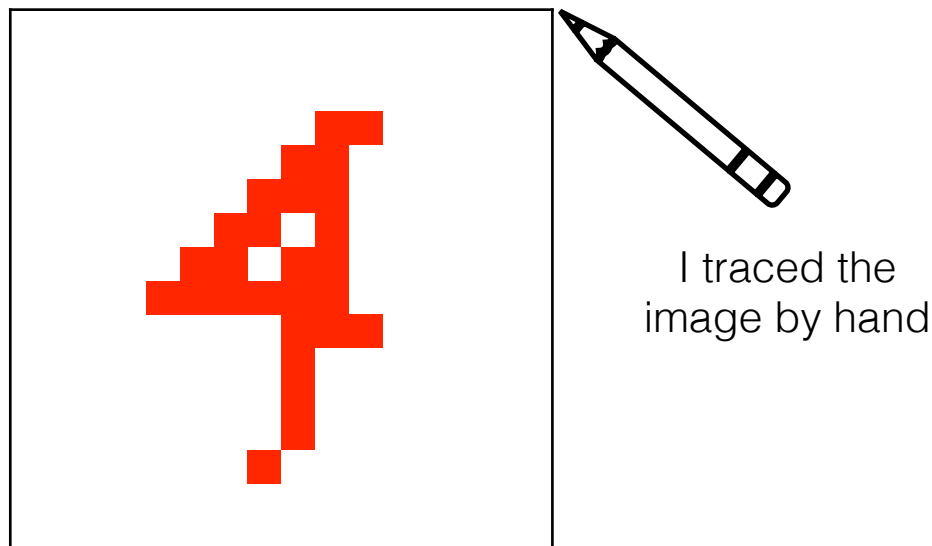
28 x 28 pixel images



Corresponding MNIST data sample

Image	Pixel 1	Pixel 2	...	Pixel 300	...	Pixel 783	Pixel 784
1	0	0	...	240	...	0	0
2	0	1	...	255	...	0	0
...	...	...	...	...	...	...	...
500k	0	0	...	231	...	0	0

16 x 16 pixel (1-bit)



$$\begin{array}{lcl} \# \text{ pixels} & 28^2 & = \\ \# \text{ bits per pixel} & 8\text{-bits} & \times \\ \# \text{ bits per image} & & \underline{6272} \\ & & 6272 \\ \# \text{ bits per image} & 16^2 \bullet 1 = & \div \\ \text{Compression factor} & & \underline{24.5} \end{array}$$

Compression isn't hard, but  
autoencoder can do better



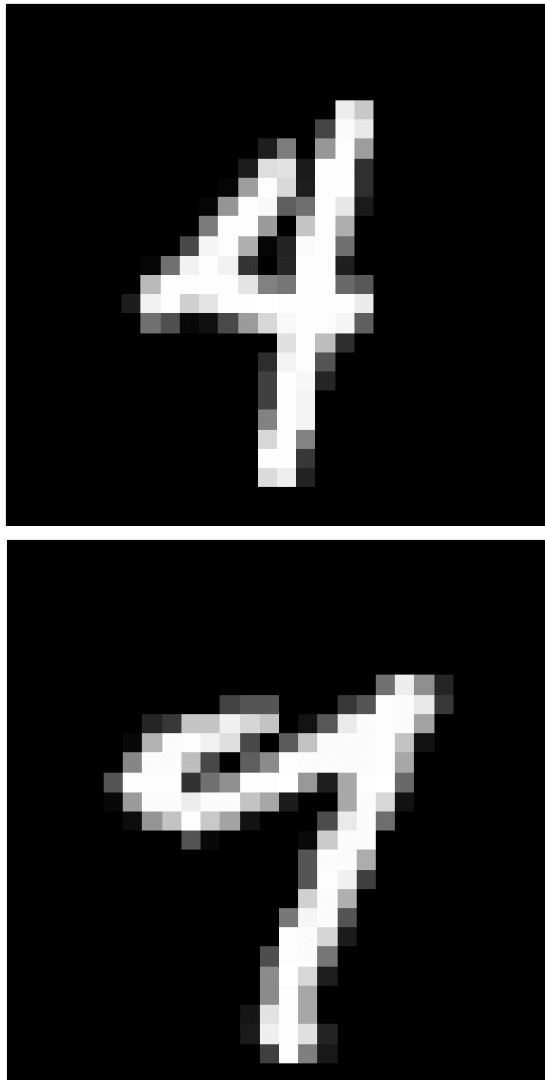
# Autoencoder recap

Example: handwritten digits

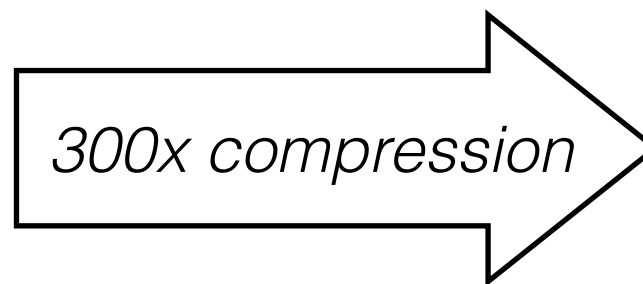
(8-bit pixels on 28x28 grid)

- Teach it 0, 1, 2, 3, 4 with a sample (doesn't know about 9!)

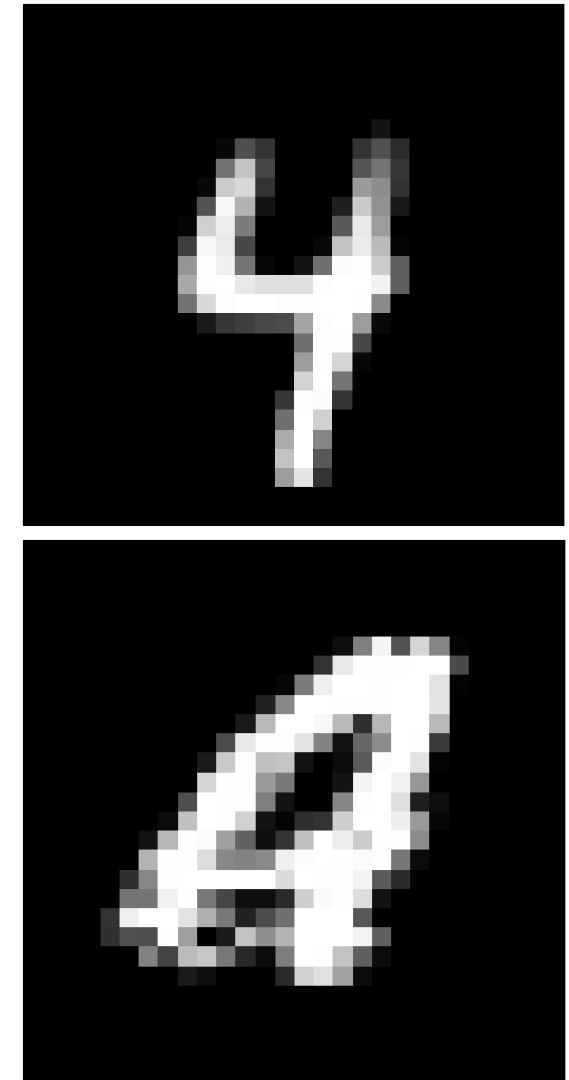
784 variables (8-bit)



1 variable (20 bit)



784 variables (8-bit)



Meaning

- Input-output distance is relatively small = good compression
- Input-output distance is relatively large = bad compression

# Strategy

How should we call the autoencoder?

## Training

done offline on CPU

Neural network-  
based training

Decision tree-  
based training

### *Anomaly detection*

- Govorkova et al.,  
Nat. Mach. Intell. 4 (2022) 154
- CMS Collaboration,  
Comput. Softw. Big Sci. 8 (2024) 11

### *Data compression*

- Guglielmo et al.,  
IEEE Trans. Nucl. Sci. 68 (2021) 2179

**Not a comprehensive list**

*Anomaly detection*  
• Gupta (for ATLAS)  
Pheno Symposium 2025  
<https://indico.global/event/812/contributions/126571>

*Data compression*  
• **This talk**

### *Anomaly detection*

- Roche, TMH et al.,  
Nat. Comm. **15** (2024) 3527  
<http://doi.org/10.1038/s41467-024-47704-8>
- Ercikti & TMH using VHDL  
**This talk**

## Deployment

for online on FPGA

Neural network-  
based design

Decision tree-  
based design

*Not sure if useful but  
certainly possible*

# I'll start with these

Training

done offline on CPU

Neural network-  
based training

Decision tree-  
based training

Deployment

for online on FPGA

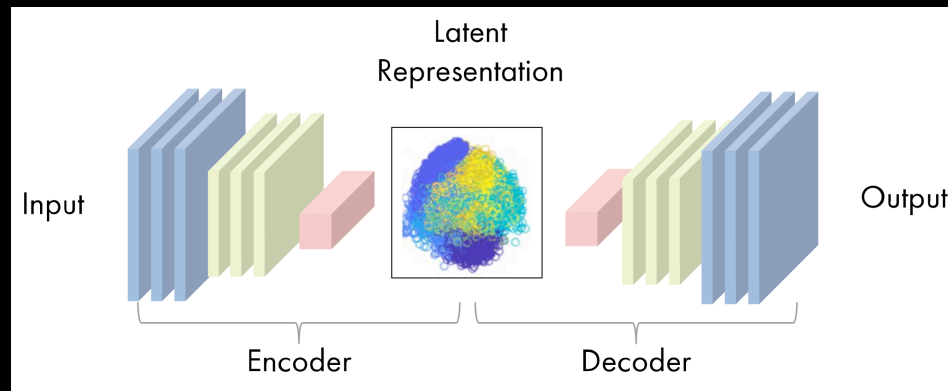
Neural network-  
based design

Decision tree-  
based design

# Autoencoders

## NN AE

- Training is a black box, done offline
- Latent space is complex



From CMS Machine Learning Group  
<https://cms-ml.github.io/documentation/training/autoencoders.html>

- FPGA version simplified for anomaly at CMS
- FPGA version can optionally skip latent sp.

From CMS Public Note, DP-2023/079  
[https://cds.cern.ch/record/2876546/files/DP2023\\_079.pdf](https://cds.cern.ch/record/2876546/files/DP2023_079.pdf)

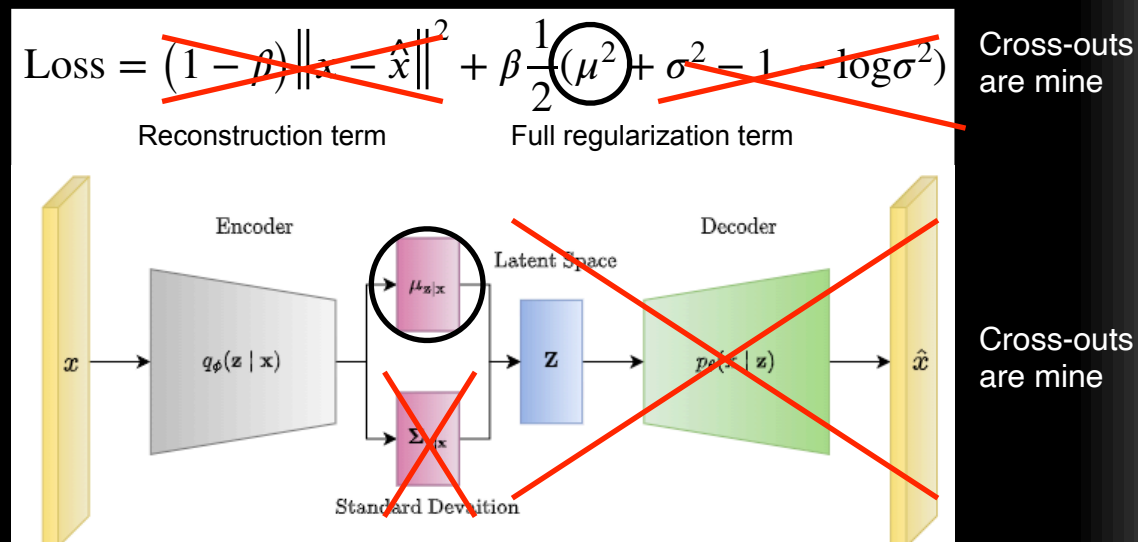
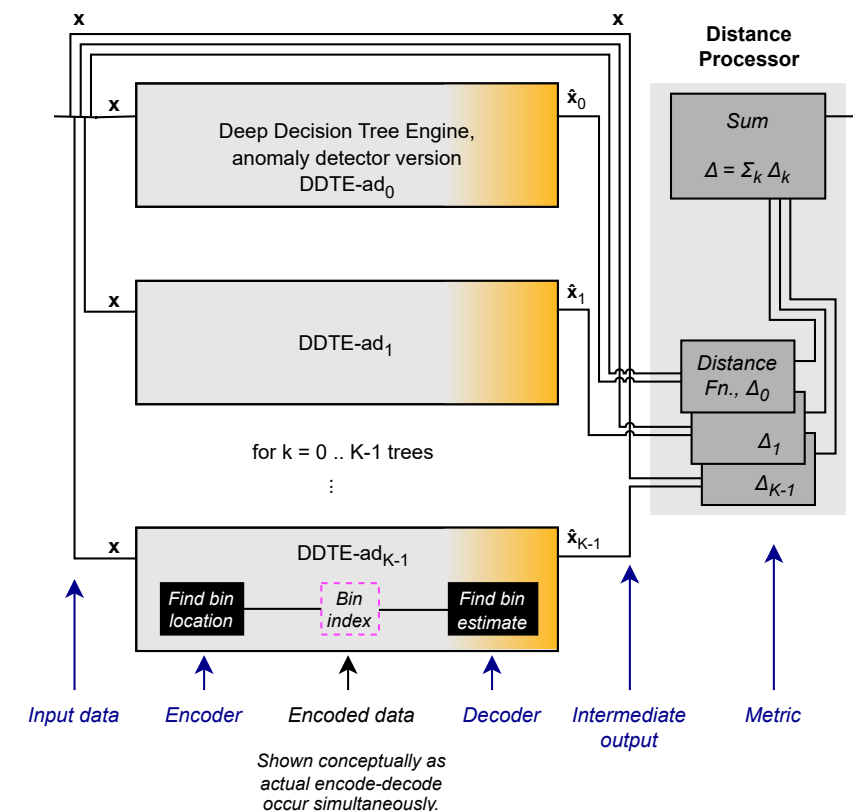
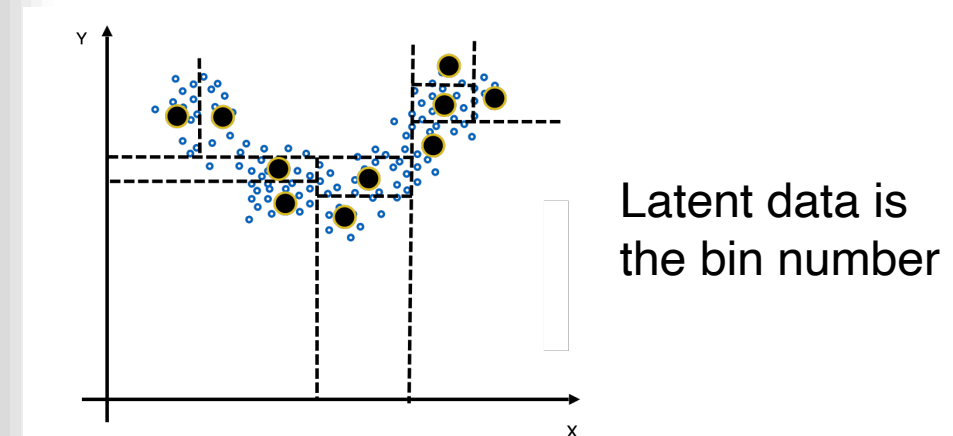


Image from  
<https://medium.com/@rushikesh.shende/autoencoders-variational-autoencoders-vae-and-beta-vae-ceba9998773d>

## "Starcoder" DT-AE

- Training is sampling of 1d pdfs
- Latent space is simple / interpretable



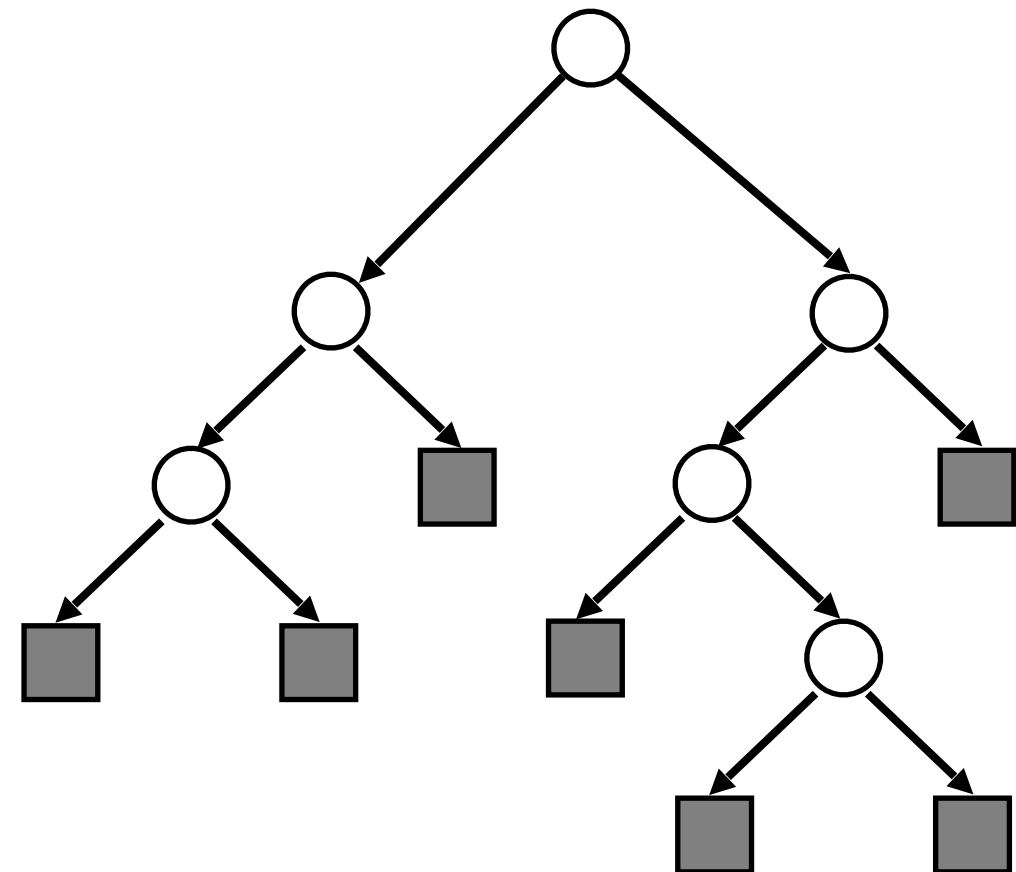
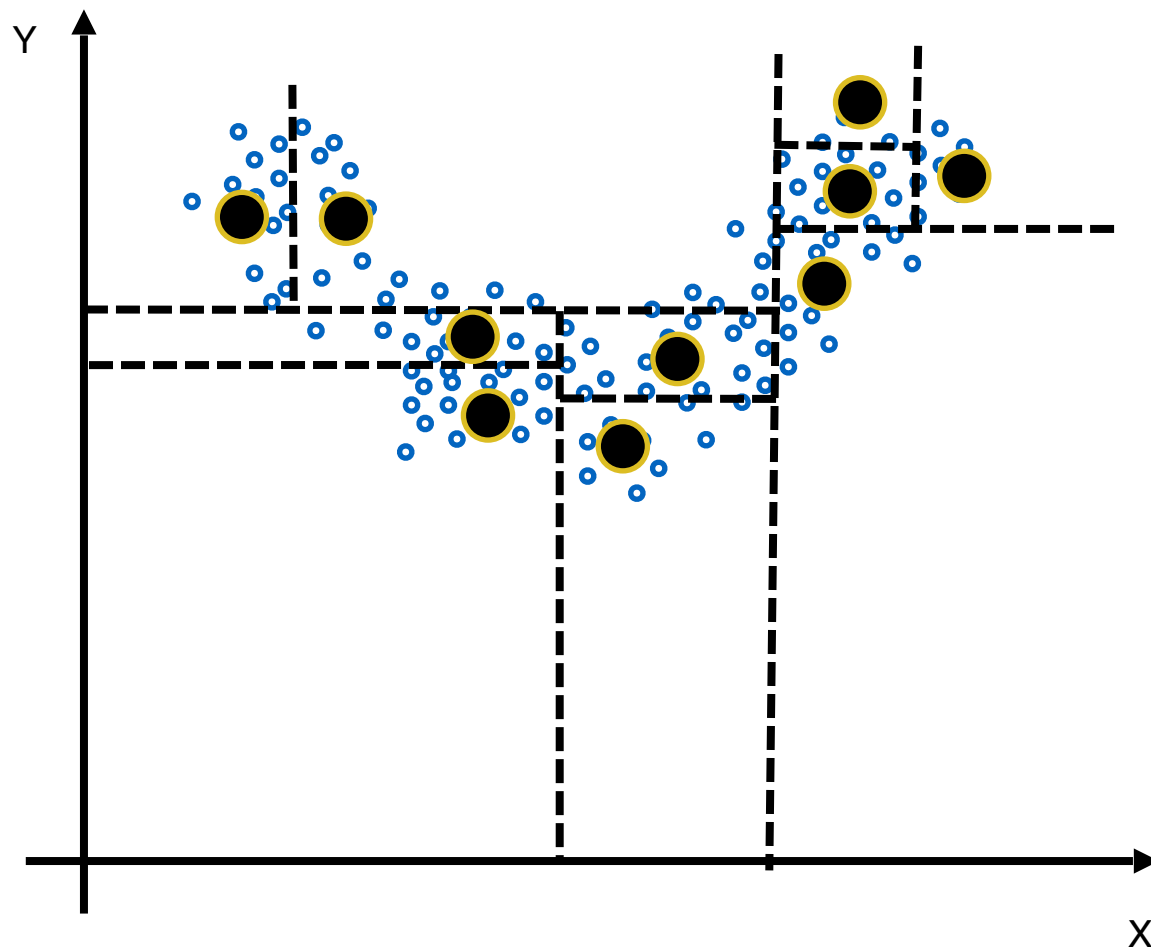
# DT-AE training developed in-house

Training by sampling 1d projection of input variables

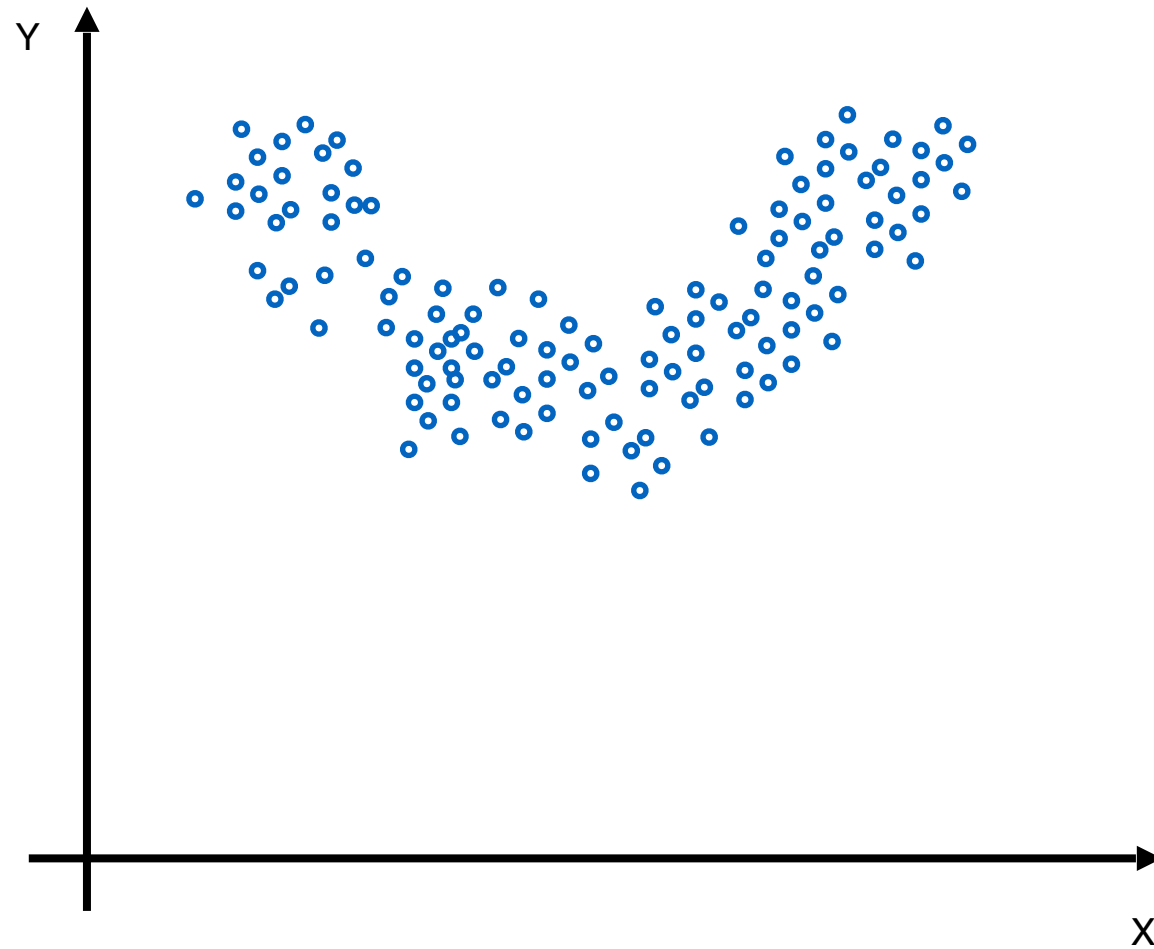
- It's doing density estimation
- Encoding: Input variables  $\rightarrow$  which bin it's in

Decoding returns “reconstruction point”

- Decoding: Bin  $\rightarrow$  median of the training data in that bin

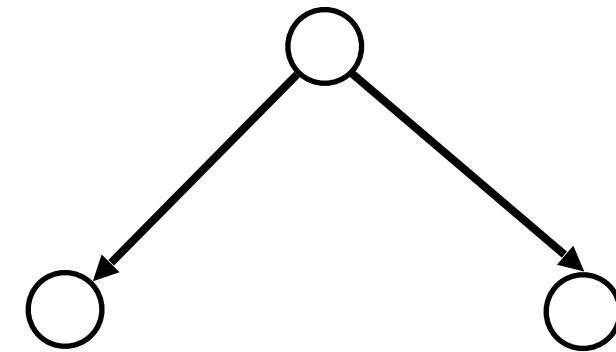
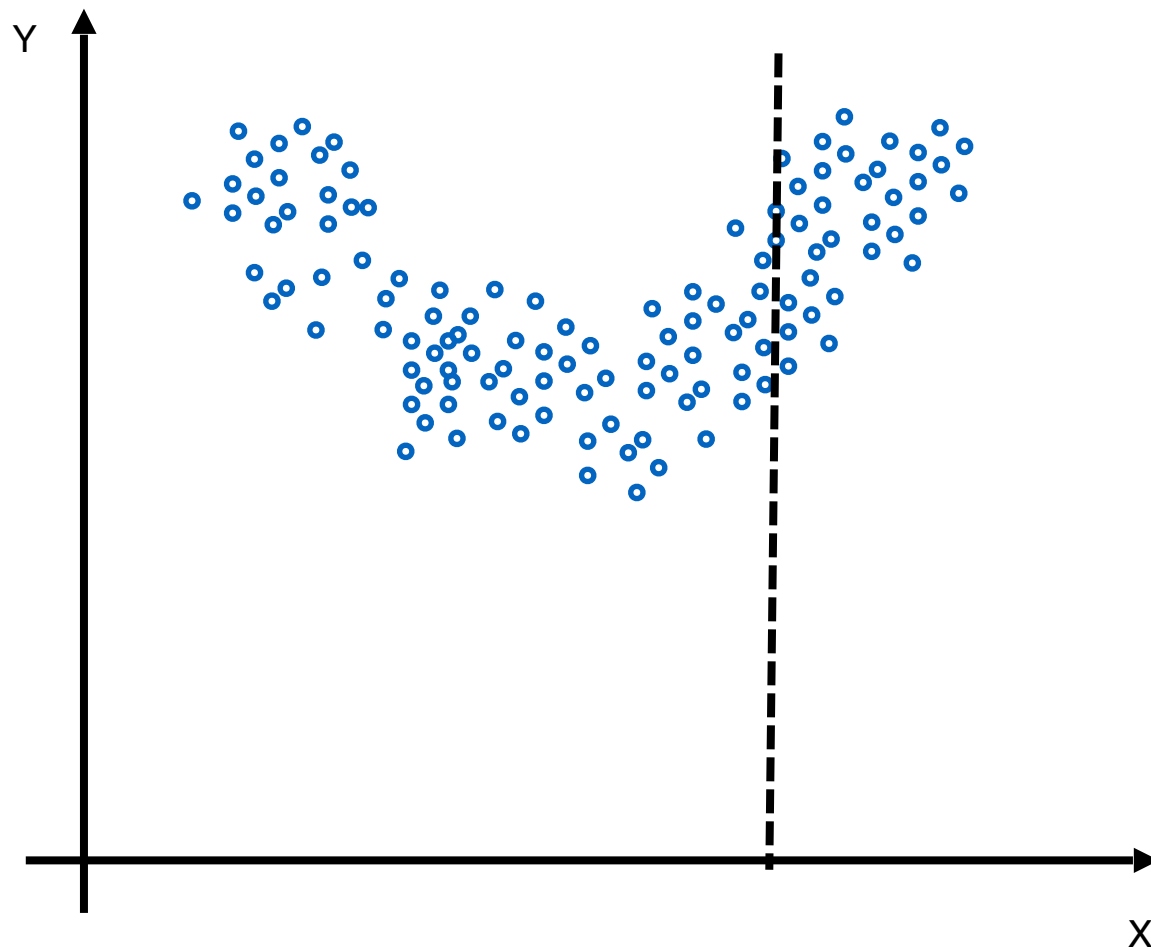


# Start over

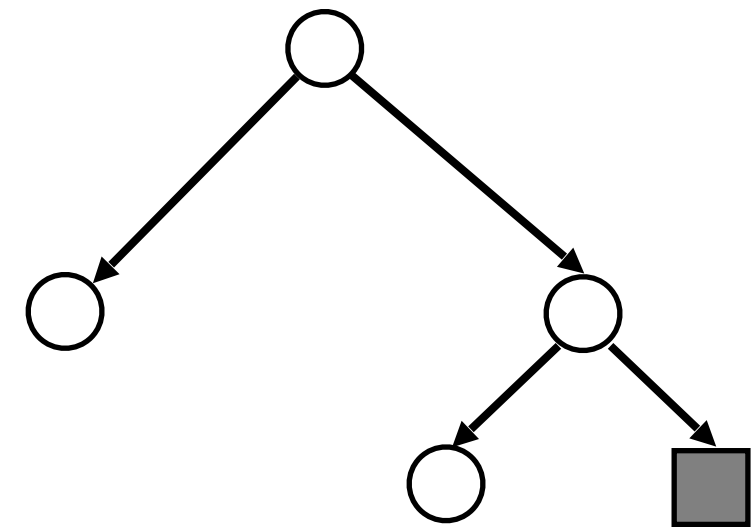
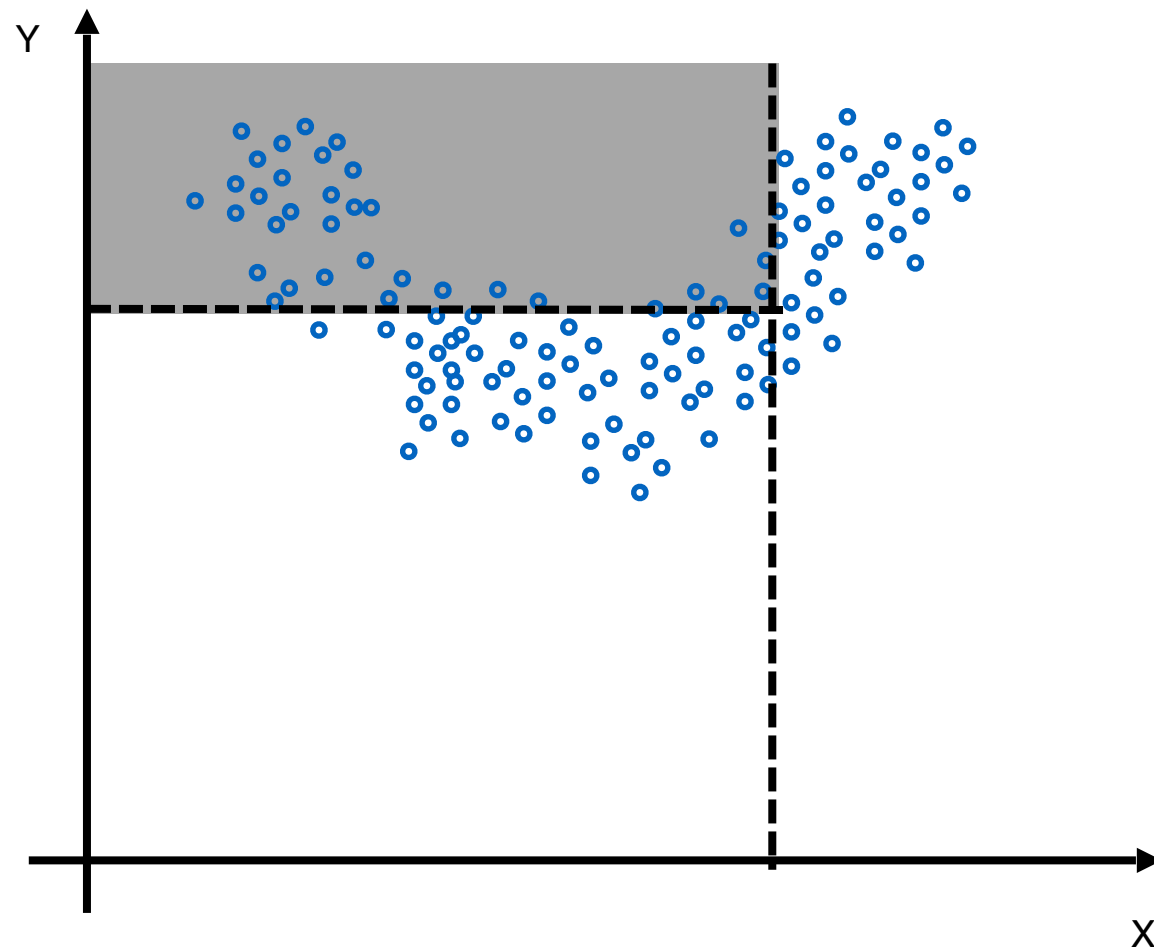




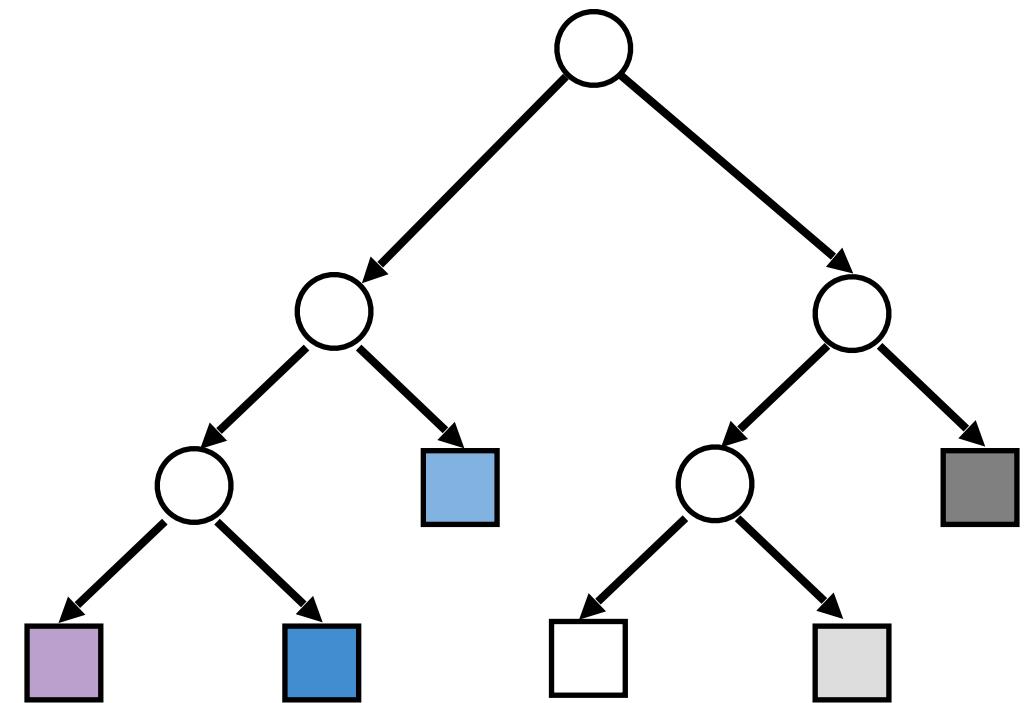
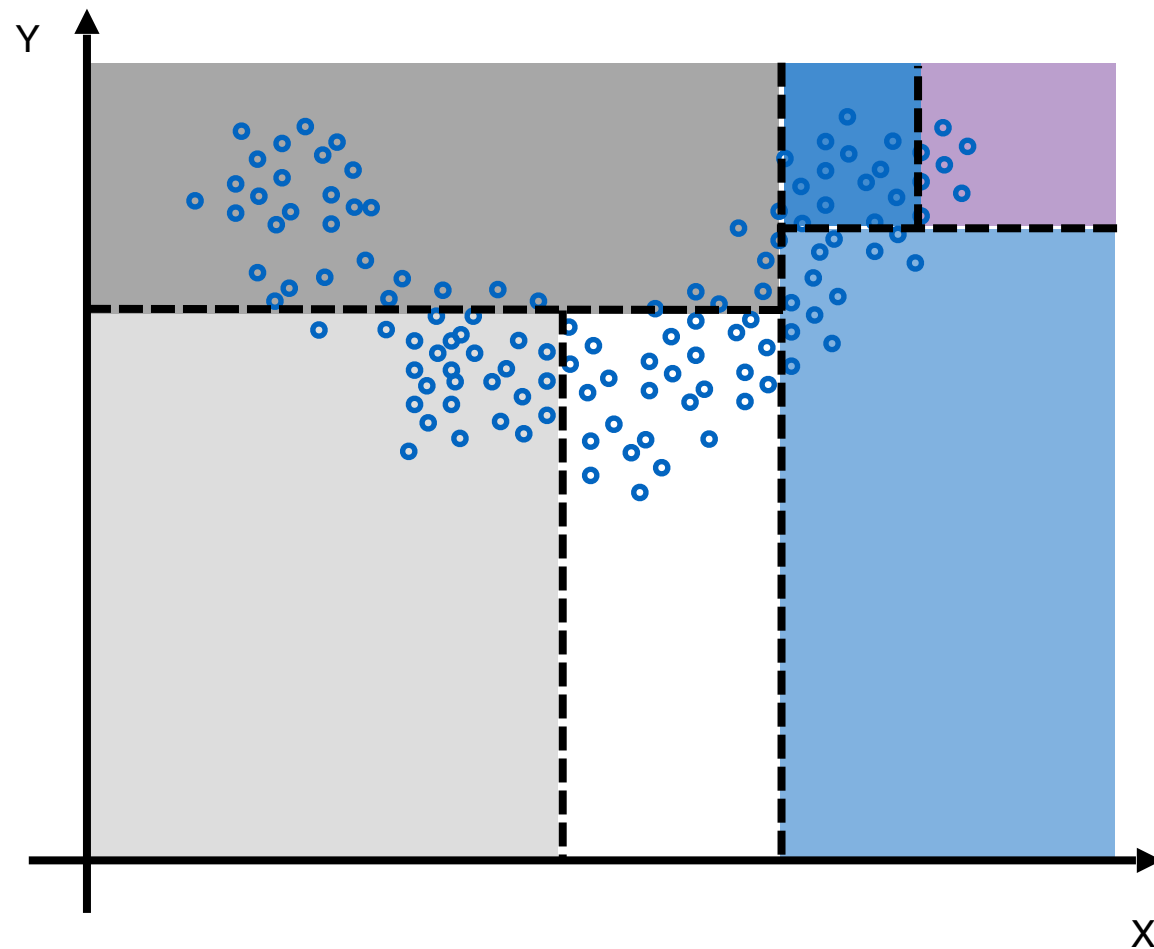
# Choose 1 variable, sample a cut



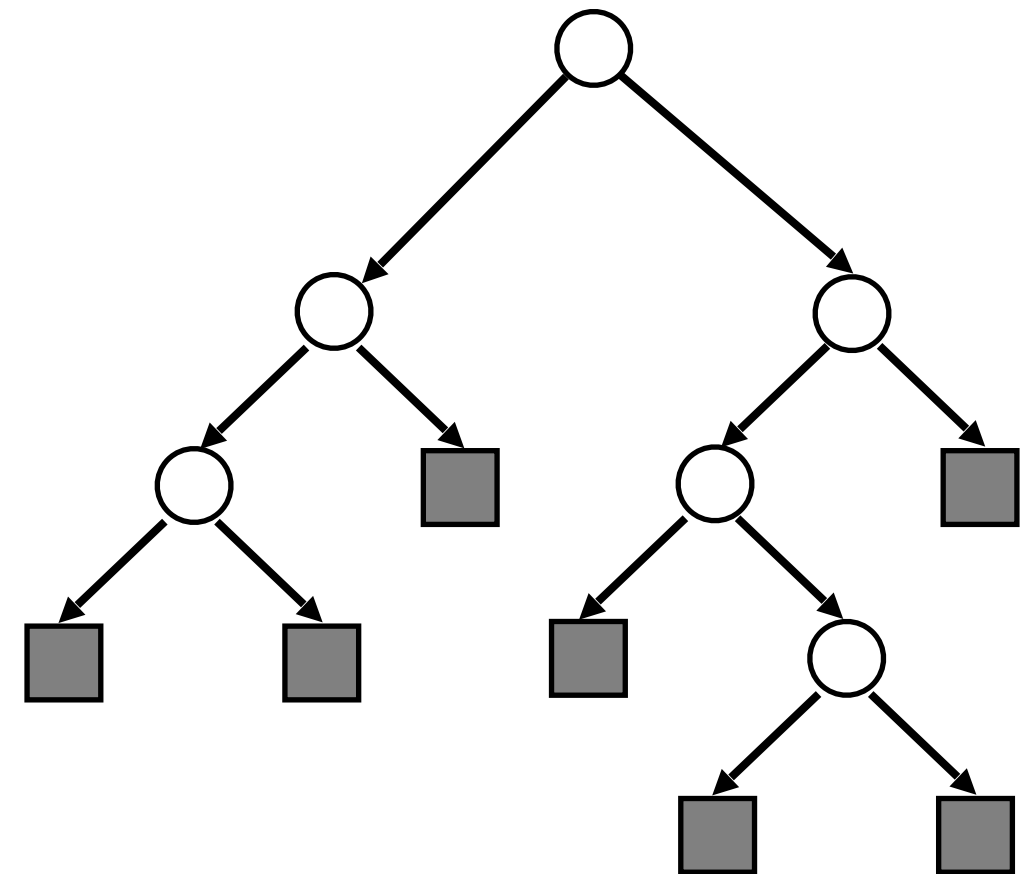
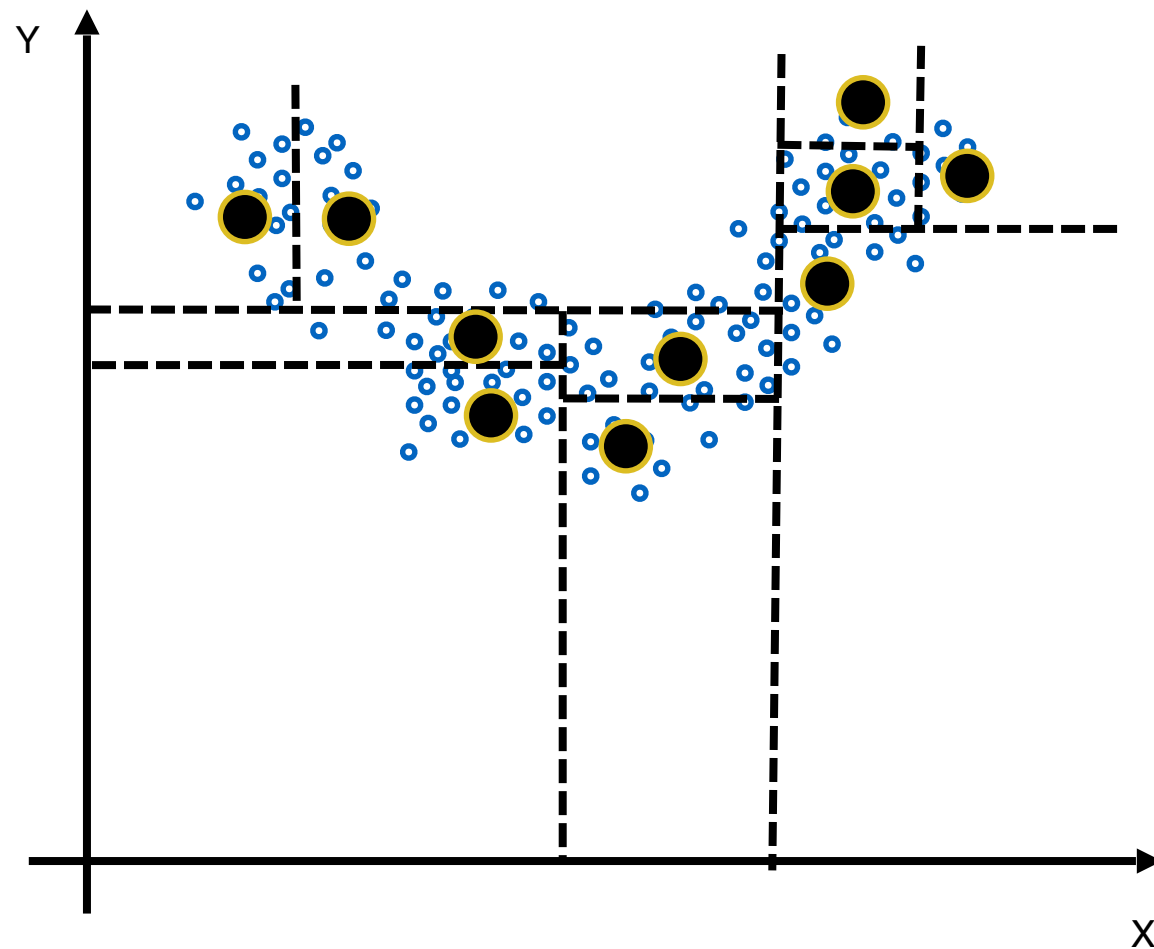
# Repeat



# Bins



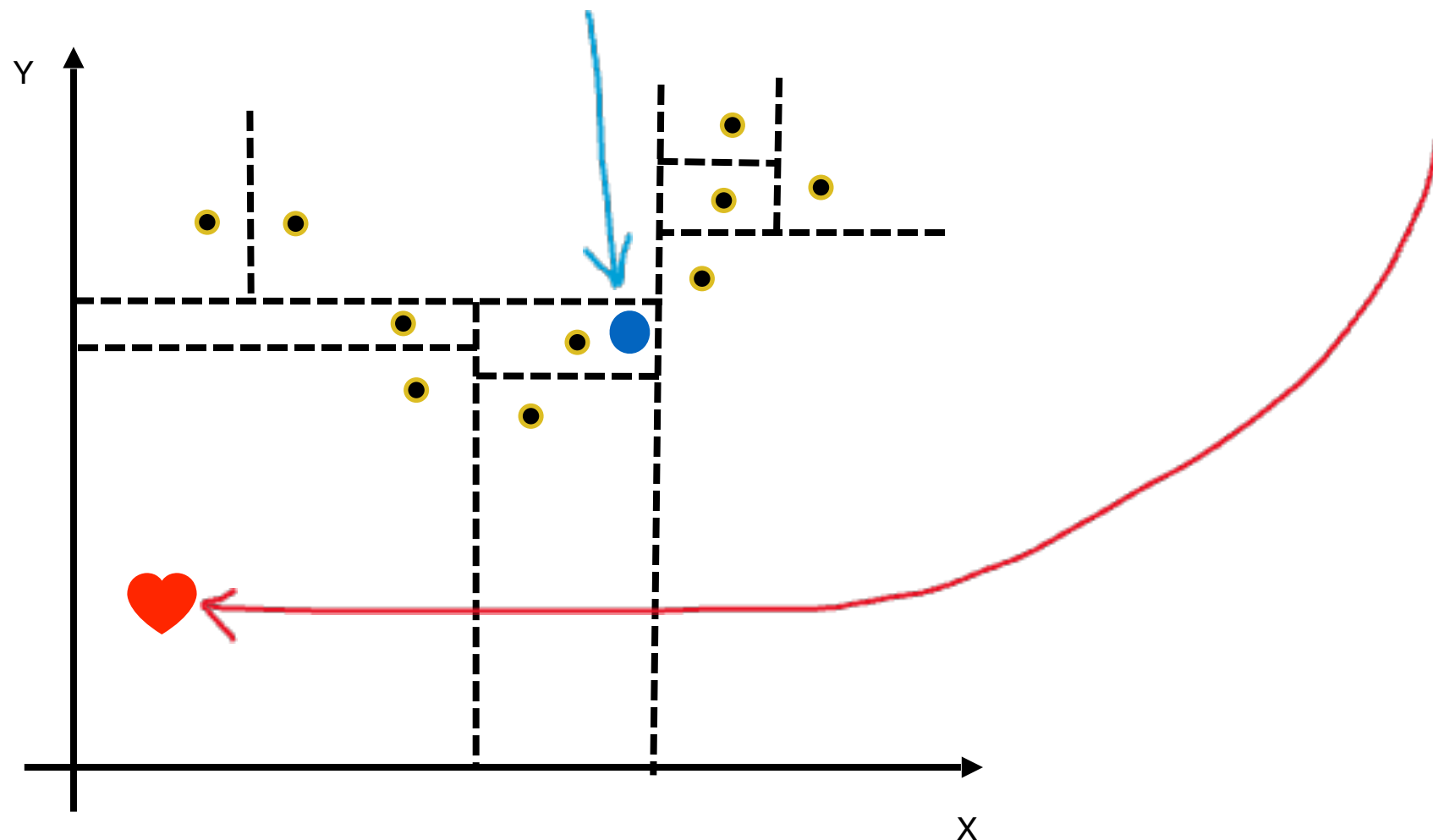
# Choose median value



# AE becomes anomaly detector

How does this detect anomalies?

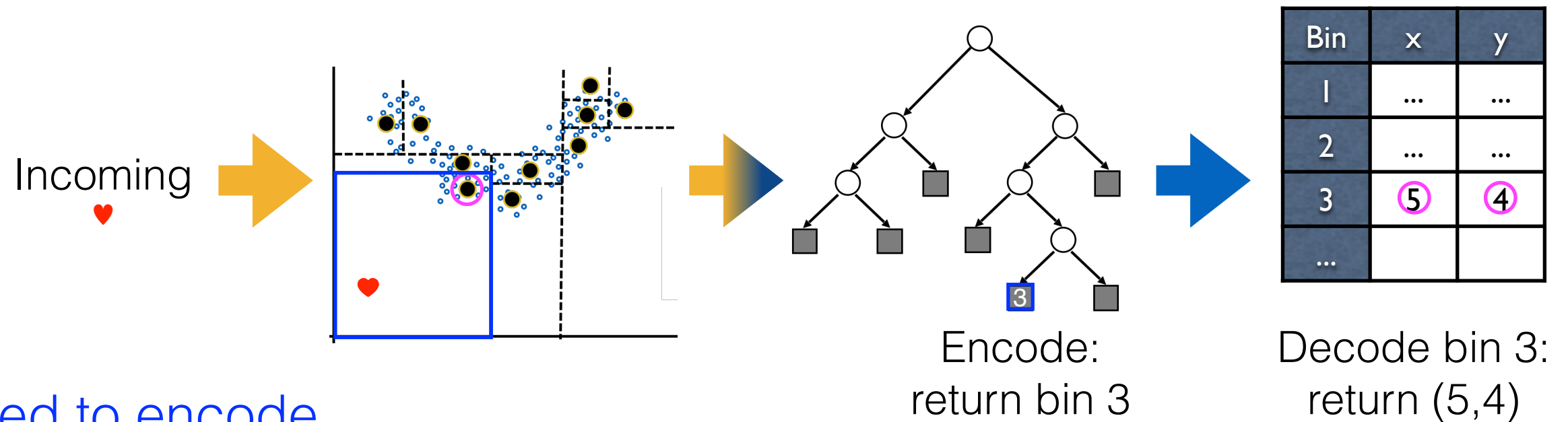
- Define: Distance between input – output = anomaly score
- Non-anomaly
  - Input is similar to training data
  - Will likely land in a small bin  $\rightarrow$  close to the reconstruction point
- Anomaly
  - Input is not similar to training data
  - Will likely land in a large bin  $\rightarrow$  far from the reconstruction point



# Realized we can skip latent space

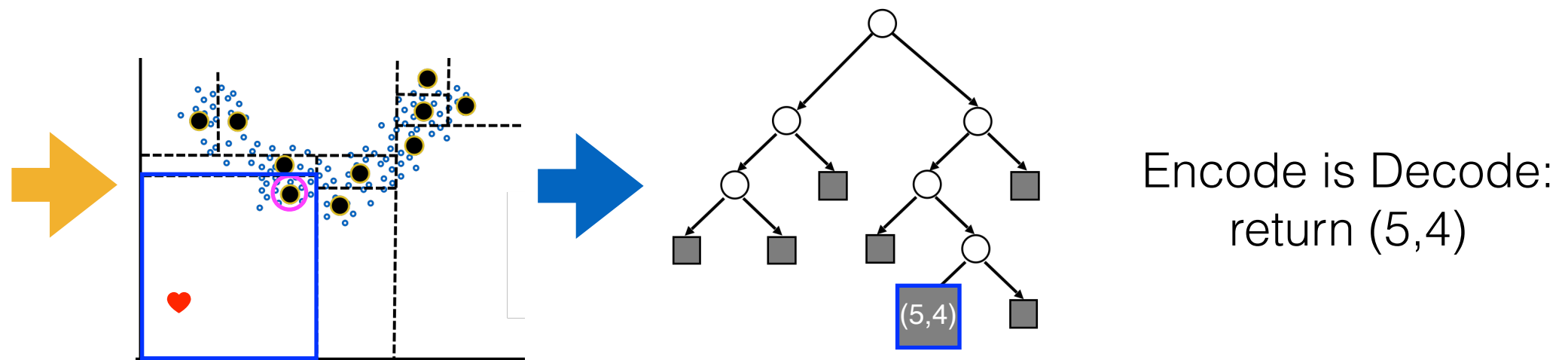
## Decode?

- Encode: input var  $\rightarrow$  bin #
- Decode: bin #  $\rightarrow$  coord.



## No need to encode

- Starcode: input var  $\rightarrow$  coord.





# Starcoder vs. hls4ml

## Works well

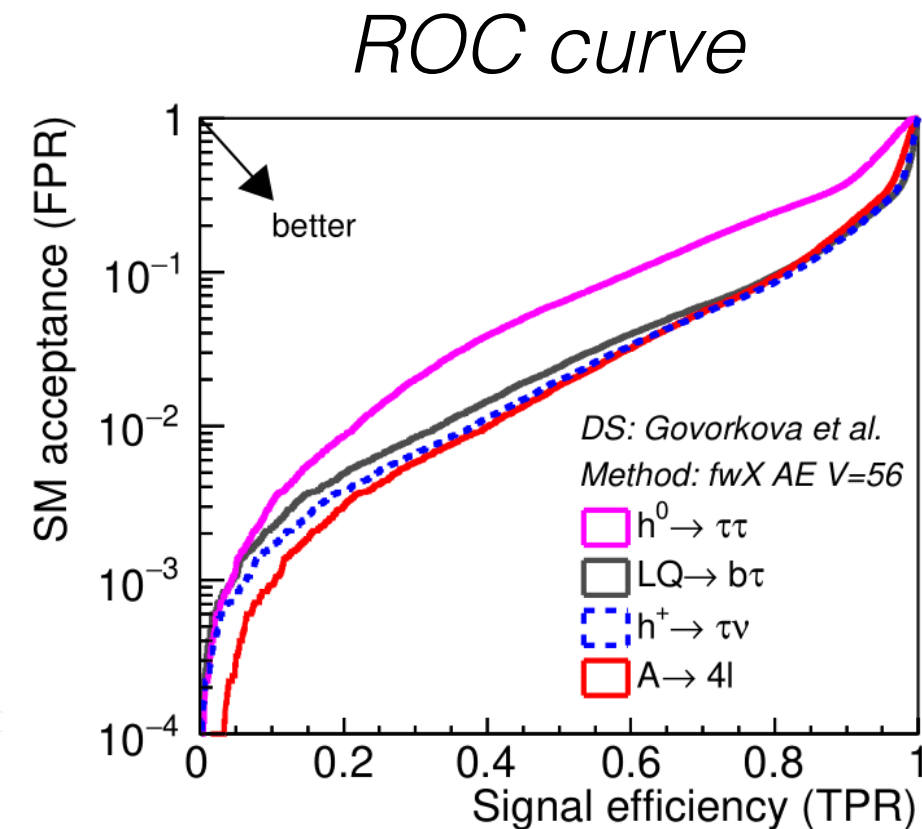
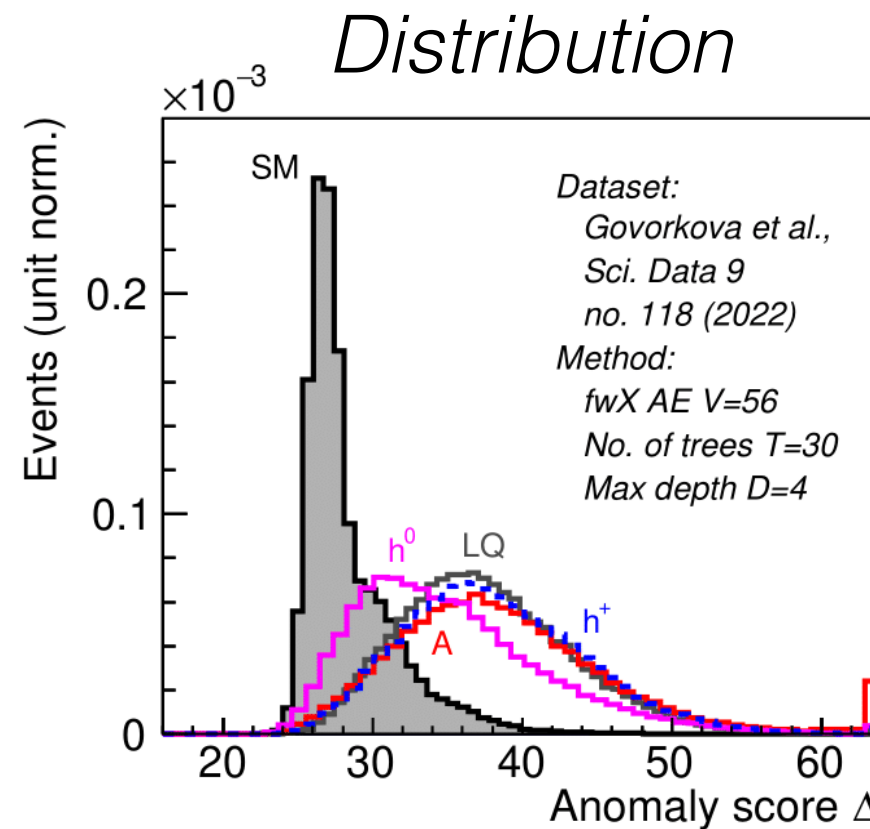
- Physics (plots)
- FPGA (table)

## Comparison

- hls4ml NN-AE  
*Nature Mach. Intell.* 4 (2022) 154
- Physics: comparable AUC
- FPGA results

## Update

- Starcoder result here uses HLS trees. Recently upgraded to VHDL trees (next slide)



	hls4ml	starcoder-hls
Clock speed	200 MHz	200 MHz
Latency	80 ns	30 ns
Interval	5 ns	5 ns
FF	0.5%	0.6%
LUT	3%	9%
DSP	1%	0.8%
BRAM	0.3%	0

# HLS vs. VHDL trees

## Starcoder got an update


- HLS High-Level Synthesis
- VHDL VHSIC Hardware Description Language  
Very High-Speed Integrated Circuit

~~C code~~ → RTL / VHDL / Verilog  
write VHDL directly

- Tree design based on  
Serhiayenka, TMH et al.  
NIM A **1072** (2025) 170209  
<http://doi.org/10.1016/j.nima.2025.170209>

## Results

- VHDL is 3-5x more efficient
- Finishing up final testbench
- Will write-up in proceedings



	starcoder-hls	starcoder-vhdl
Clock speed	200 MHz	320 MHz
Latency	30 ns	25 ns
FF	0.6%	3.2x smaller
LUT	9%	5.3x smaller
DSP	0.8%	0
BRAM	0	0

# Now the diagonal - 1

Training

done offline on CPU

Neural network-  
based training

Decision tree-  
based training

*Anomaly detection*

- Gupta (for ATLAS)

Pheno Symposium 2025

<https://indico.global/event/812/contributions/126571>

Deployment

for online on FPGA

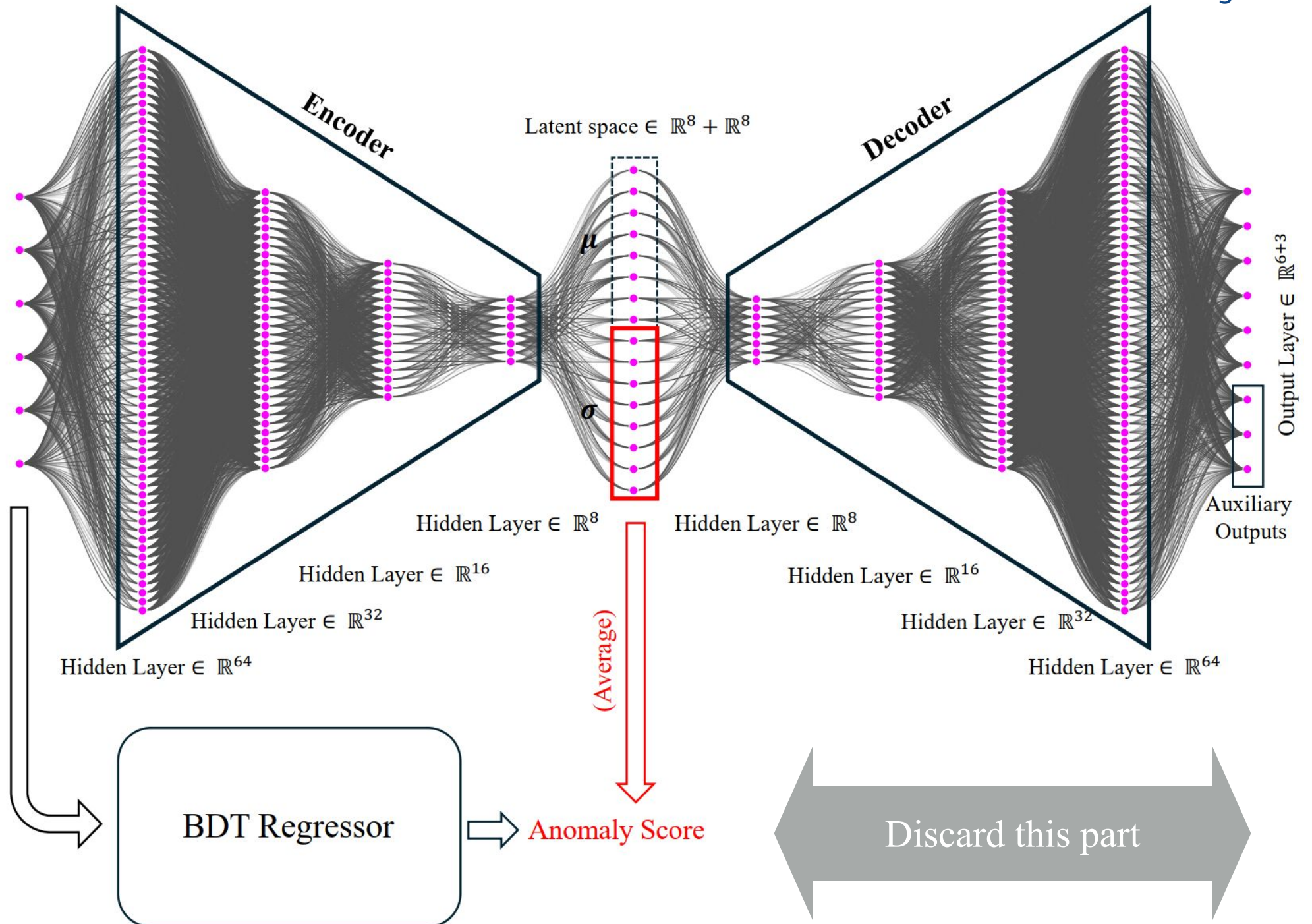
Neural network-  
based design

Decision tree-  
based design

# VAE variational autoencoder



Train VAE



Regress

- Carlson, TMH et al.  
JINST **17**, P09039 (2022)  
<http://doi.org/10.1088/1748-0221/17/09/P09039>

Deploy Trees on FPGA



# How it looks like in an experiment

I'm not representing ATLAS today

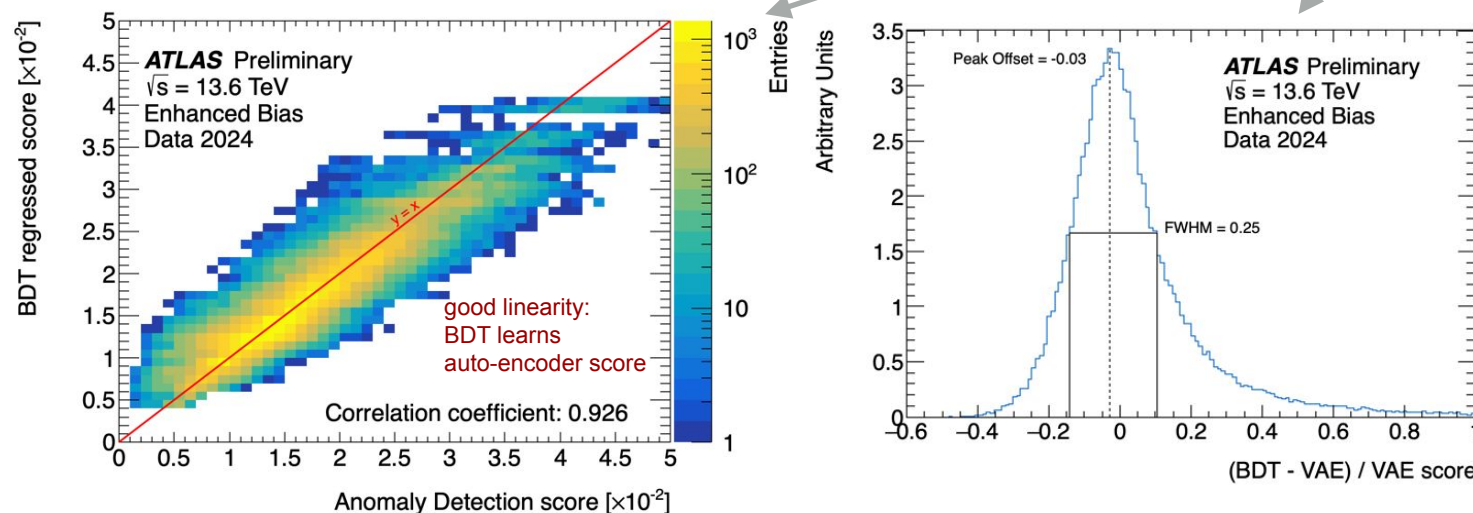
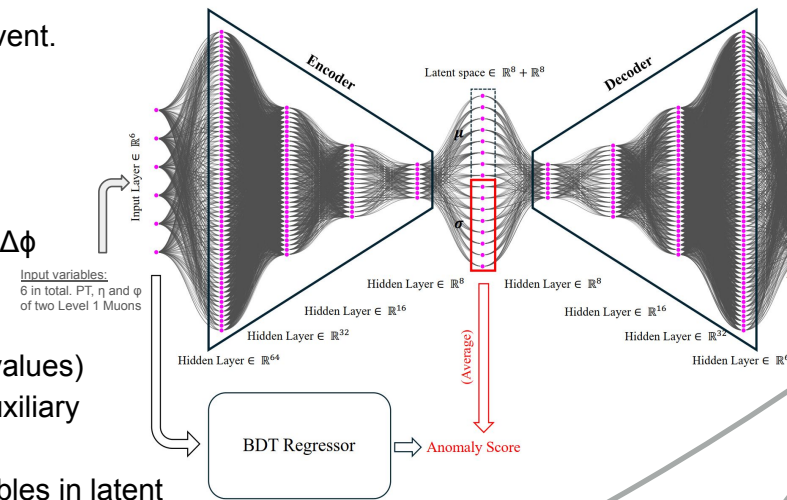
## Training: Variational Autoencoder

R. Gupta



8

- Data: enhanced bias 2024 data ([ATL-DAQ-PUB-2016-002](#))
- Preselection: 2MU3VF
- At most 3 leading muons (leading in  $p_T$ ) per event.
- Train on unlabeled dimuon events ( $p_T$ ,  $\eta$ ,  $\phi$  of muon pairs) (details in backup).
- Consider all combinations (e.g. 3 muon pairs) for one event.
- Regressed variables in the model: mass,  $\Delta R$ ,  $\Delta\phi$
- VAE architecture:
  - Encoder: 4 layers
  - Latent space: 8D Gaussian ( $\mu + \sigma \rightarrow 16$  values)
  - Decoder reconstructs inputs + predicts auxiliary vars ( $m_{\mu\mu}$ ,  $\Delta R$ ,  $|\Delta\phi|$ )
- Anomaly score = average value over 8 $\sigma$  variables in latent space.
- train a BDT regression to learn the NN score
- Same 6 variables as inputs for the BDT (2 muons'  $P_T$ ,  $\eta$ ,  $\phi$ )
- BDTs trained with TMVA: 200 trees, Depth of 20.



Source: [https://twiki.cern.ch/twiki/bin/view/AtlasPublic/MuonTriggerPublicResults#NoMAD\\_Nanosecond\\_Anomaly\\_detection](https://twiki.cern.ch/twiki/bin/view/AtlasPublic/MuonTriggerPublicResults#NoMAD_Nanosecond_Anomaly_detection)

Being commissioned at ATLAS

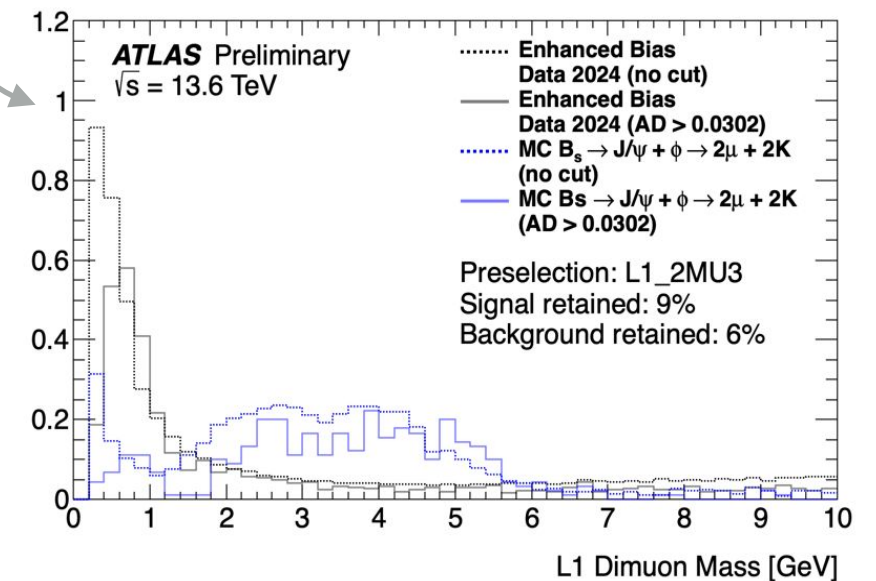
- Train VAE for 3 muons
- Executes in 1 clock tick (25 ns)

Method

- Chop-off the decoder
- Regress the latent space variables

Physics result

- Unique B physics signal at L1



- Gupta
- Pheno Symposium 2025 in Pittsburgh
- <https://indico.global/event/812/contributions/126571>

# Now the diagonal - 2

Training

done offline on CPU

Neural network-  
based training

Decision tree-  
based training

*Data compression*  
• This talk

Deployment

for online on FPGA

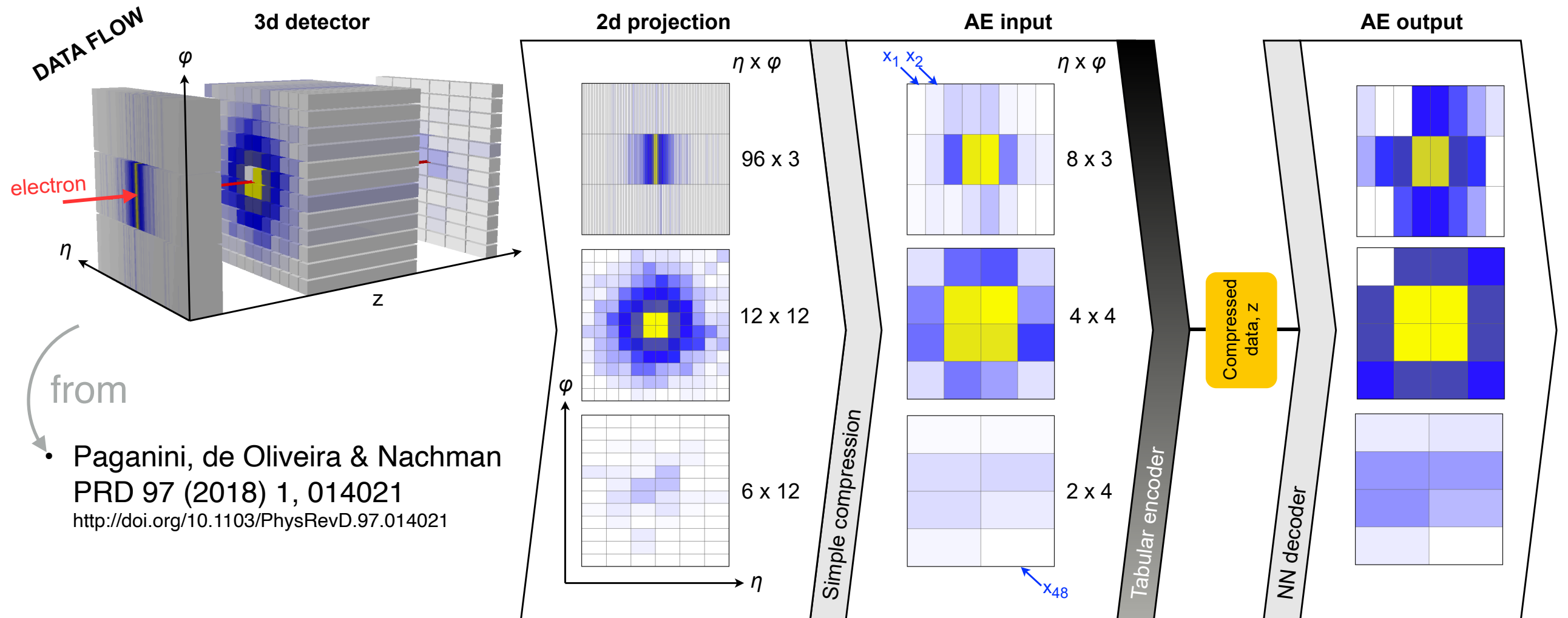
Neural network-  
based design

Decision tree-  
based design

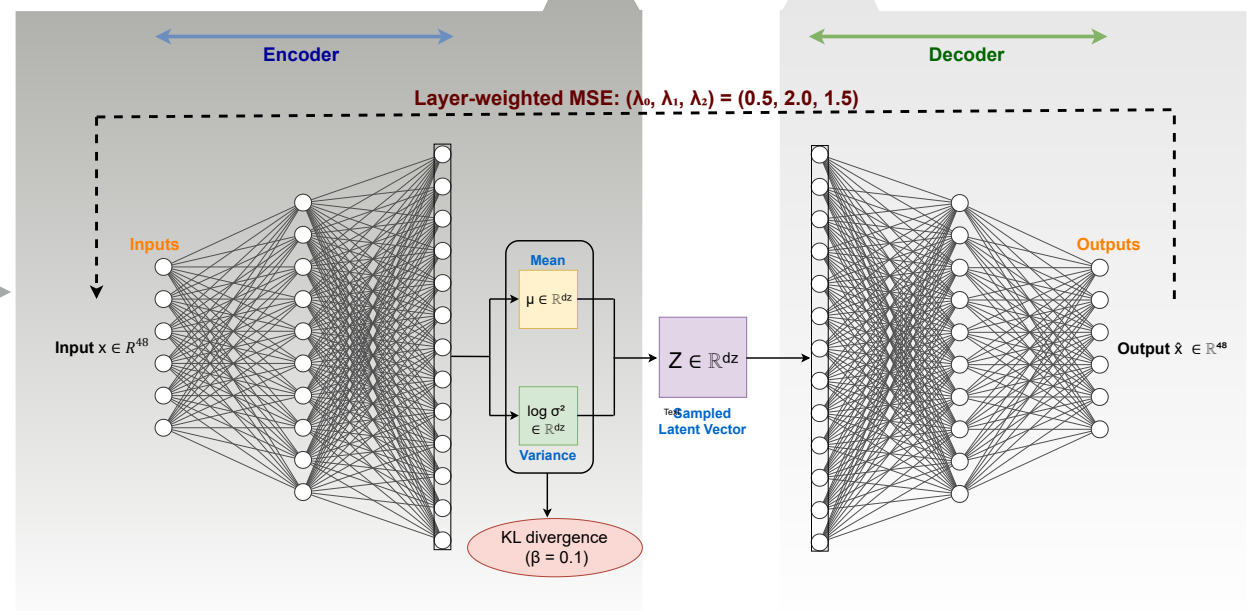


# Data compression using VAE

Use toy calorimeter dataset

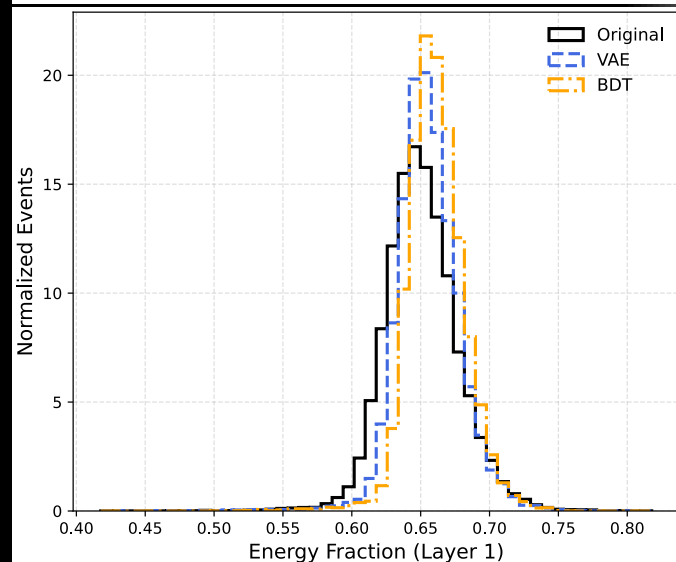
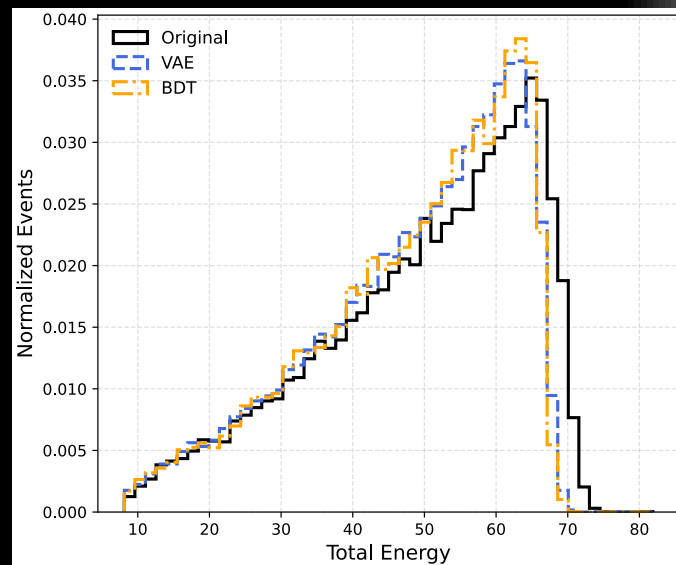


We convert this left  
part into decision trees



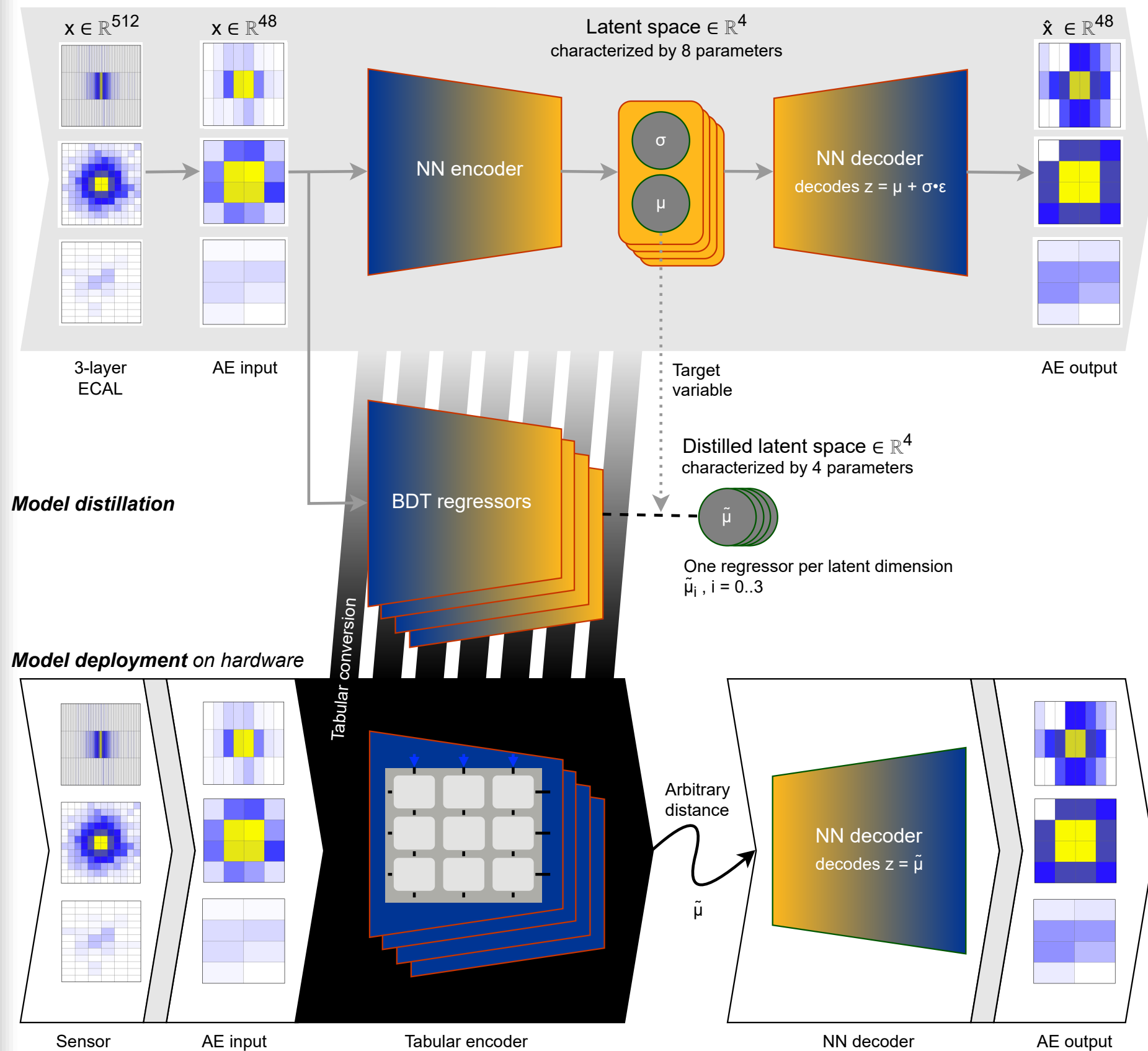
# Tabular encoder

- Physics is preserved



- Convert VAE into tabular format
- Tabular? Next slide

## Model training using simulated data

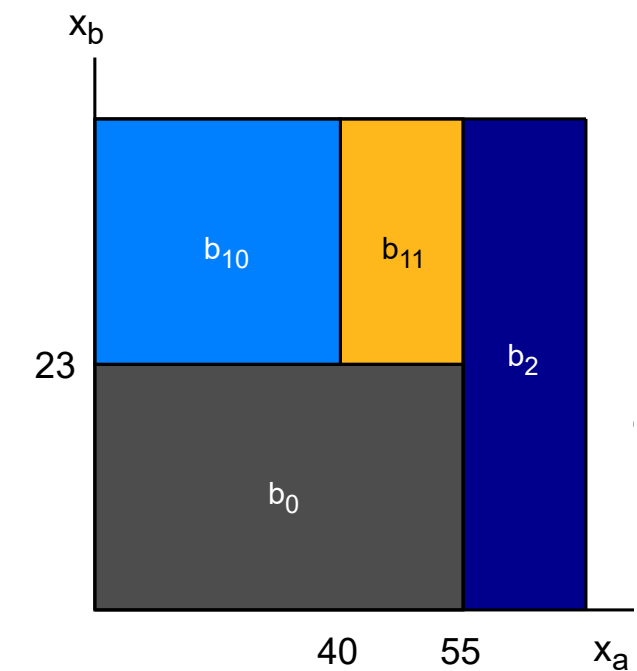
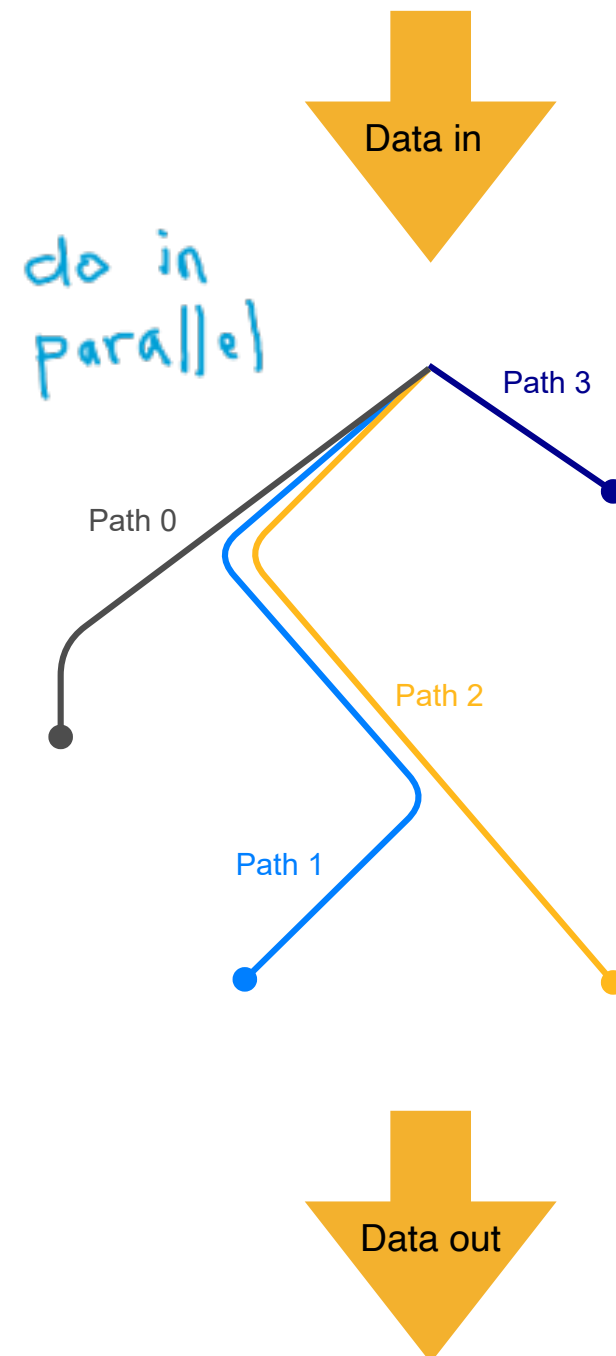
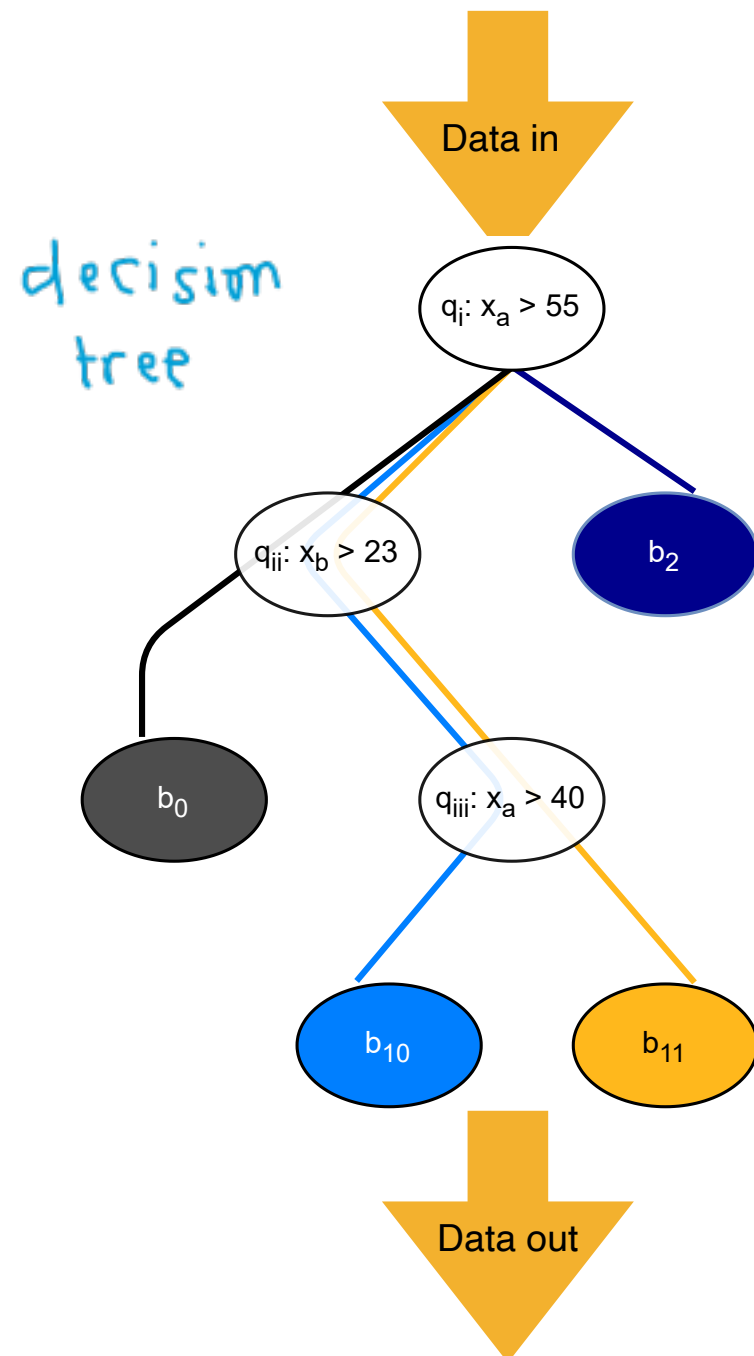


- Gupta, TMH et al., Paper in preparation

# Tabular?

## Parallel implementation of decision tree on FPGA

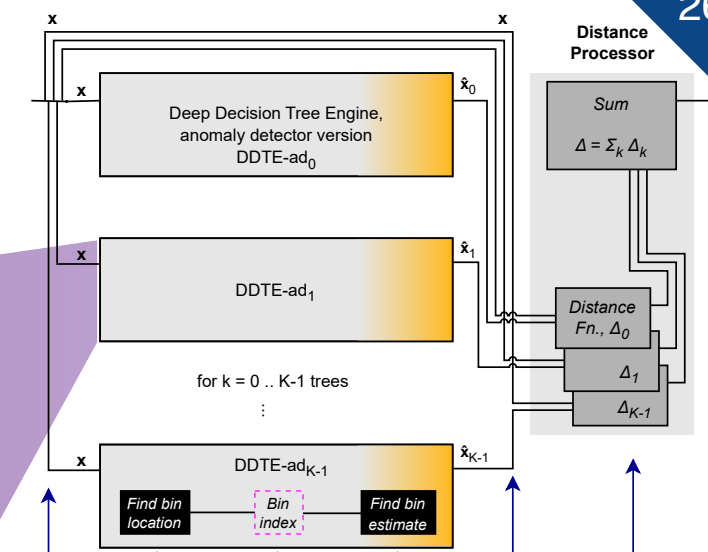
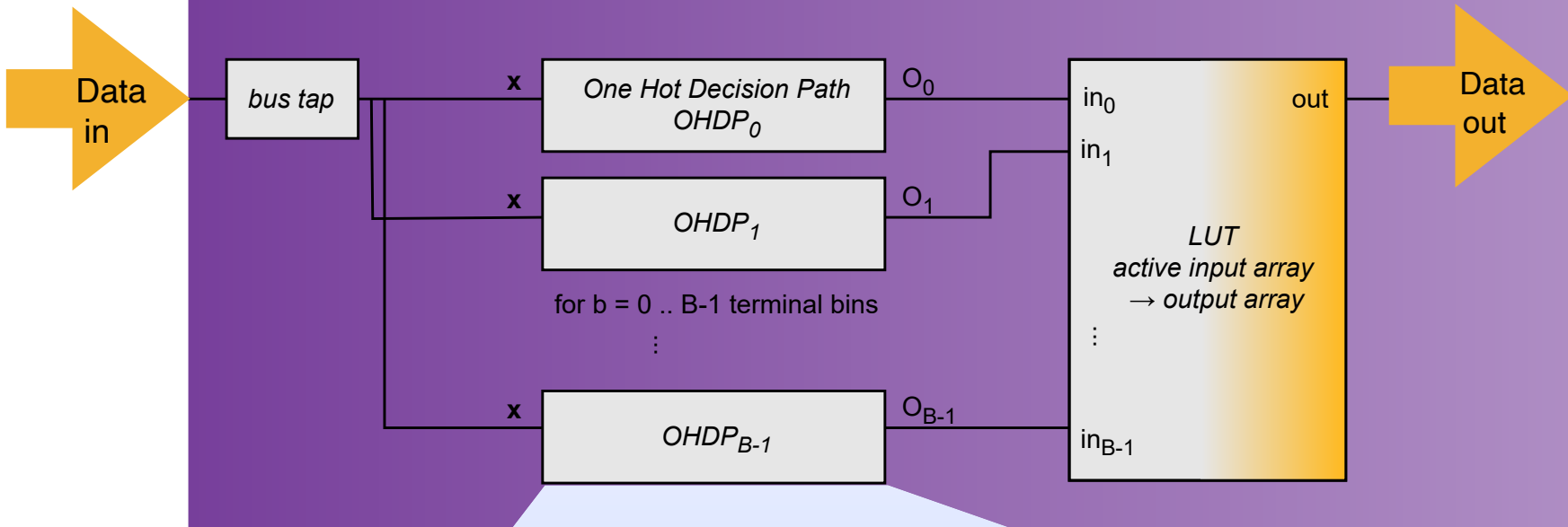
- B. Carlson, TMH et al.  
J. Instrum. 17 (2022) P09039  
<http://doi.org/10.1088/1748-0221/17/09/P09039>



Destination bin	Decision path
b <sub>0</sub>	not(q <sub>i</sub> ) and not(q <sub>ii</sub> )
b <sub>2</sub>	q <sub>i</sub>
b <sub>10</sub>	not(q <sub>i</sub> ) and q <sub>ii</sub> and not(q <sub>iii</sub> )
b <sub>11</sub>	not(q <sub>i</sub> ) and q <sub>ii</sub> and q <sub>iii</sub>

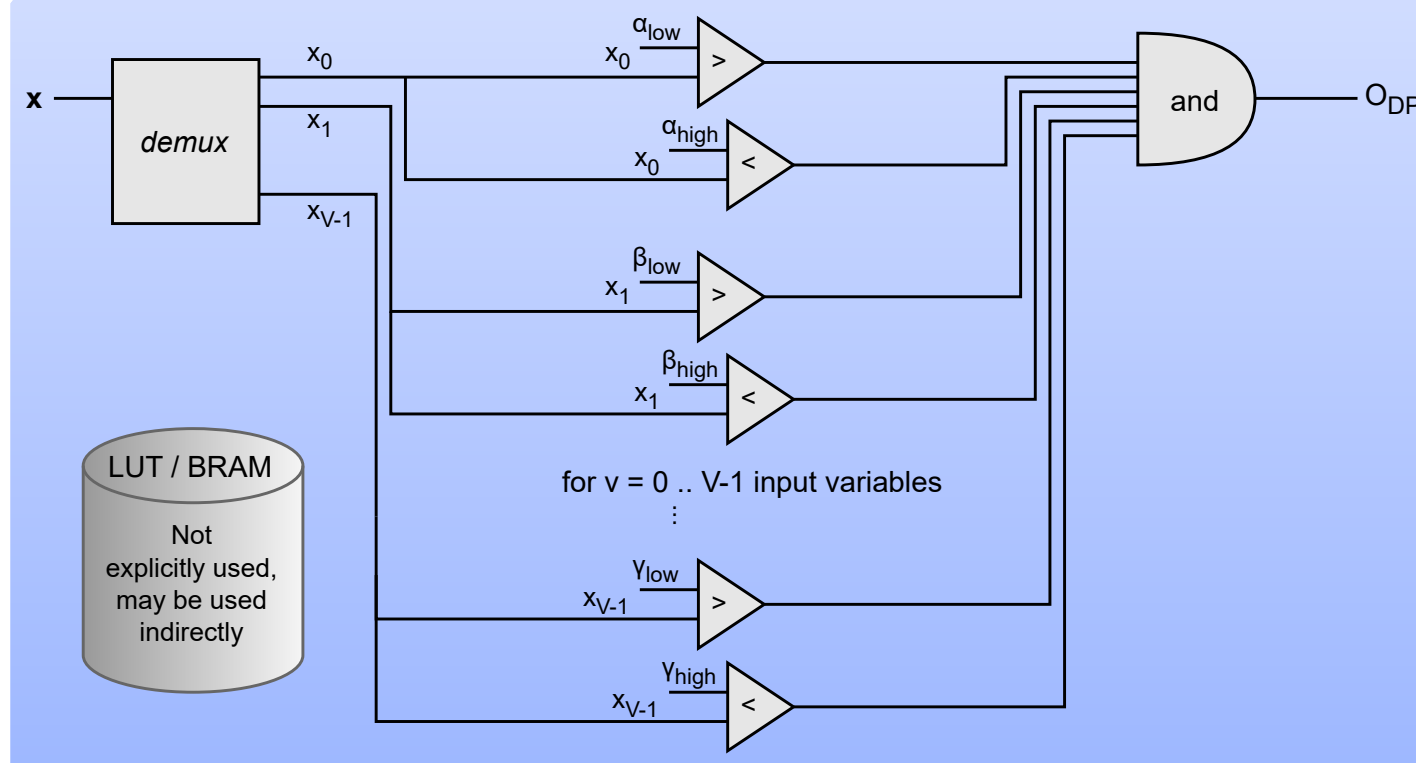
# Tabular design

## Parallel implementation



Using Xilinx Ultrascale+ VU9P (vcu118)  
at 200 MHz

Feature	Value
Latency	2 clock ticks (50 ns)
Interval	1 clock tick (25 ns)
Flip-flops (FF)	10399 (0.44%)
Look-up tables (LUT)	13274 (1.1%)
Digital signal processors (DSPs)	0
Block-RAM (BRAM)	9 (0.36%)
Ultra-RAM (URAM)	0



Serhiayenka, TMH et al.  
NIM A **1072** (2025) 170209  
<https://doi.org/10.1016/j.nima.2025.170209>

# Fun slide - question

Can you guess what these are? Hint: Belgian traffic signs

<https://btsd.ethz.ch/shareddata>

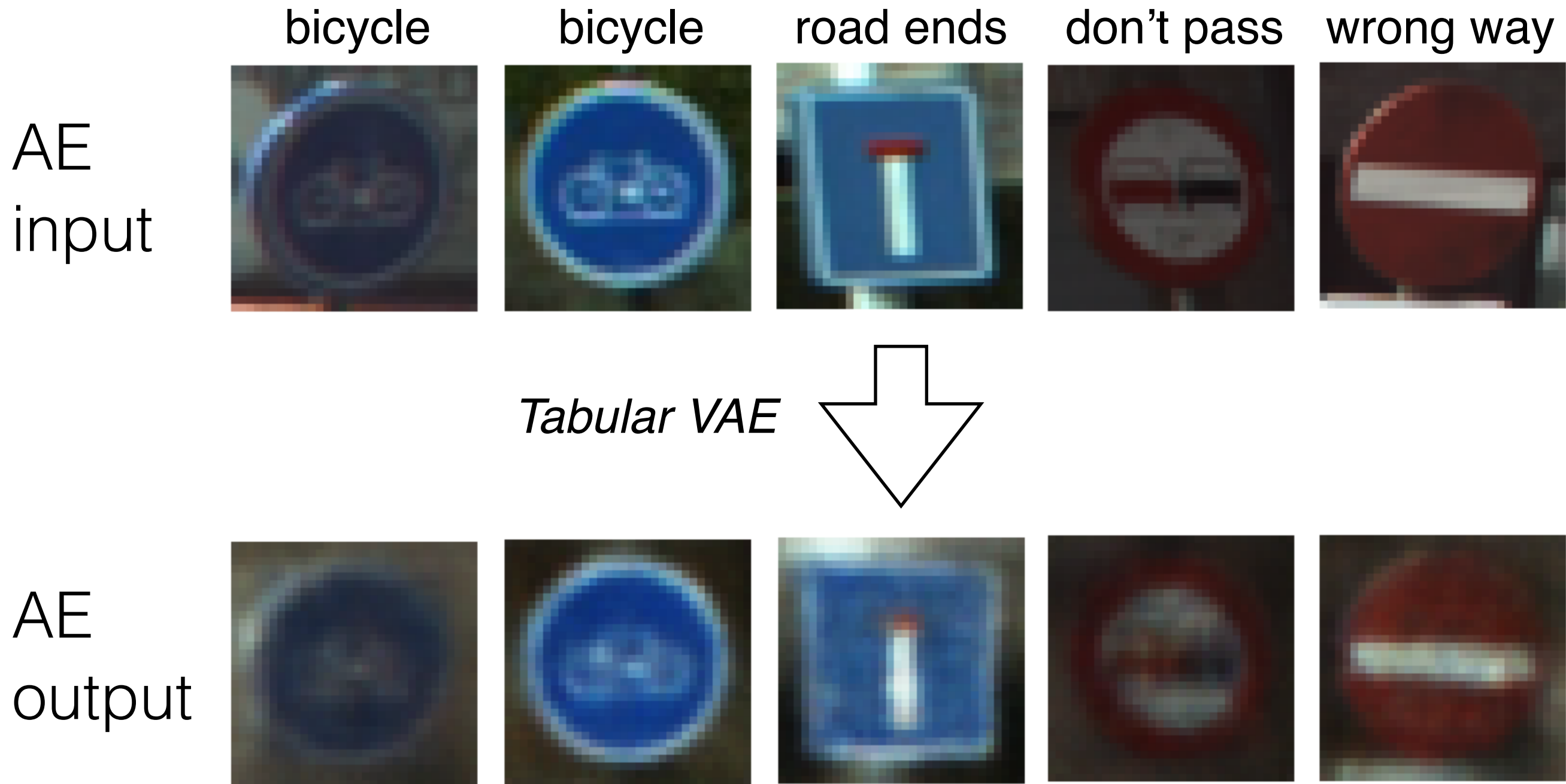
AE  
output



- Timofte et al.,  
IEEE Workshop on Appl. of Comp. Vision, WACV 2009  
<https://btsd.ethz.ch/shareddata/publications/Timofte-WACV-2009.pdf>

# Fun slide - answer

Were you right?



- Timofte et al.,  
IEEE Workshop on Appl. of Comp. Vision, WACV 2009  
<https://btsd.ethz.ch/shareddata/publications/Timofte-WACV-2009.pdf>



# Python-based code

## Availability

- [gitlab.com/PittHongGroup/fwX](https://gitlab.com/PittHongGroup/fwX)

parallel cuts (paper 1)

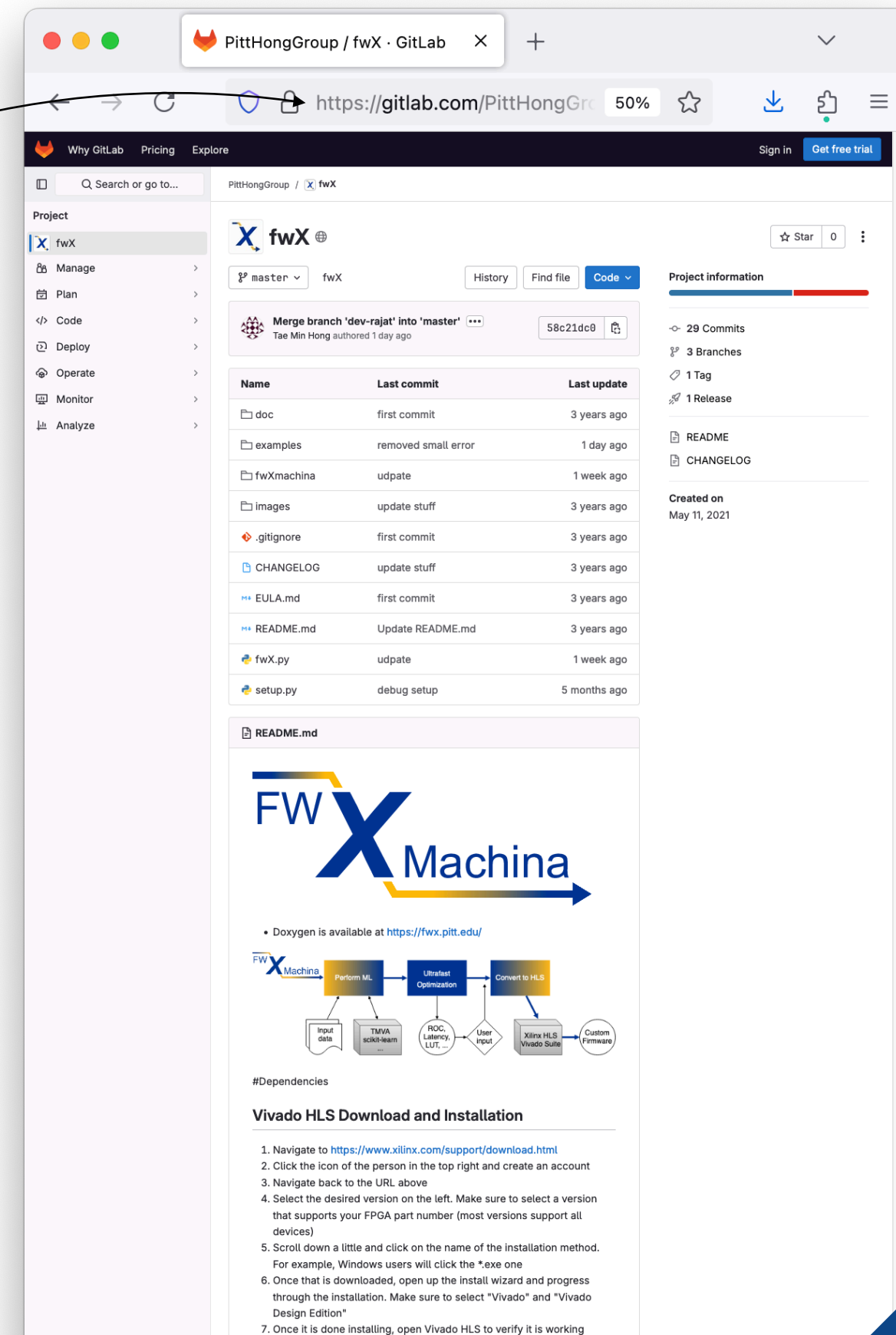
- Shared by email request

parallel paths (paper 2)

autoencoder (paper 3)

hardware tree (paper 4)

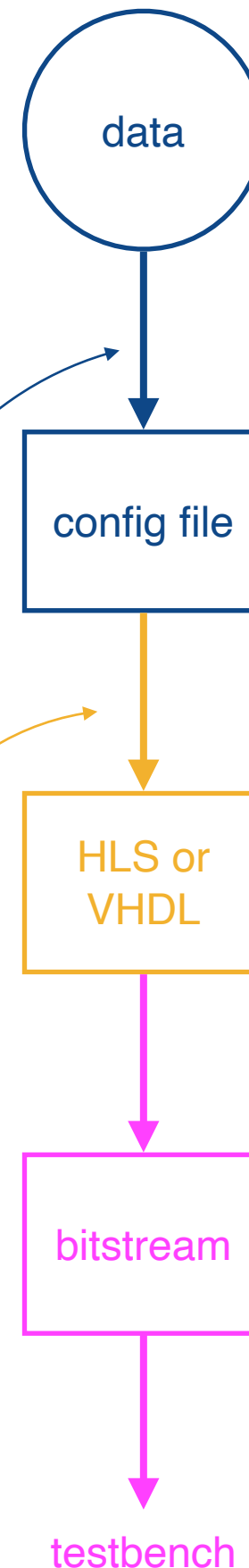
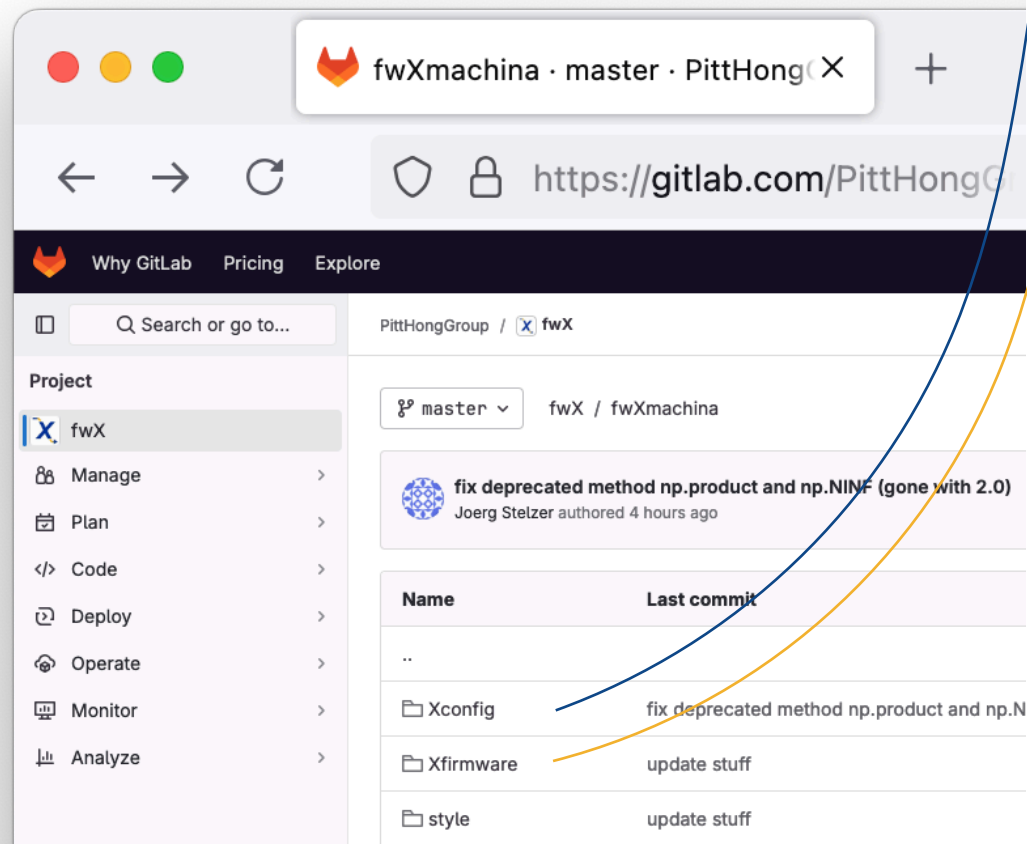
- **Collaborators welcome**



# Git structure

Same structure for all methods

- [gitlab.com/PittHongGroup/fwX](https://gitlab.com/PittHongGroup/fwX)  
parallel cuts (paper 1) - tutorial today
- Available by request  
parallel paths (paper 2)  
autoencoder (paper 3)  
hardware tree (paper 4)



- Xconfig  
creates model configuration  
**tutorial - part 1**
- Xfirmware  
writes HLS or VHDL  
**tutorial - part 2**
- Vivado  
synthesize & testbench  
**tutorial - part 3**

# FW testbench w/ IP available

<http://d-scholarship.pitt.edu/45784/>

## Screenshots in the document



# Autoencoder Firmware Testbench Tutorial

Please download Vivado 2019.2 at the following link, if you do not currently have it:

<https://www.xilinx.com/support/download/index.html/content/xilinx/en/downloadNav/vivado-design-tools/archive.html>

## Before Beginning

Before beginning, please make sure that you have (and know the location of) the autoencoder IP folder, and the VHDL testbench files:

Name	Date modified	Type	Size
 autoencoder8var_ip	2/7/2024 1:30 PM	File folder	
 tb_vhd_files	2/8/2024 11:50 AM	File folder	

## Creating New Project in Vivado

Open Vivado 2019.2 and select “create new Project.” On the following pop-up, select “next,” and you will be prompted to name the project. Name the project as you wish and choose a location to store it. Keep clicking next until you reach a page that prompts you to select the part/ board. For this tutorial, we will be using the Virtex UltraScale+ VCU118 board. After you have selected your part or board, keeping clicking “next” until you have reached the end of the setup page.

