# FW**X**MACHINA

## Nanosecond machine learning with boosted decision trees for high energy physics

University of Pittsburgh

Tae Min Hong

- *Paper* [2104.03408]
- *Info*    fwx.pitt.edu
- *Code*   gitlab.com/PittHongGroup/fwX

Pheno 2021

May 24, 2021

https://indico.cern.ch/event/982783/sessions/396894/

# Nanosecond machine learning event classification with boosted decision trees in FPGA for high energy physics

T.M. Hong,* B.T. Carlson, B.R. Eubanks, S.T. Racz,[†] S.T. Roche, J. Stelzer, and D.C. Stumpp[‡]

Department of Physics and Astronomy
University of Pittsburgh

*Undergraduate researchers*

May 17, 2021

† Received his BS last month!
‡ Now a Pitt PhD student in EE

- Intro

- Algorithm structure

- Firmware design

- Physics results (simulated)

_____

- Backup slides

- Intro
  - Machine learning at Level-1 trigger

- Algorithm structure
  - Bit integer representation
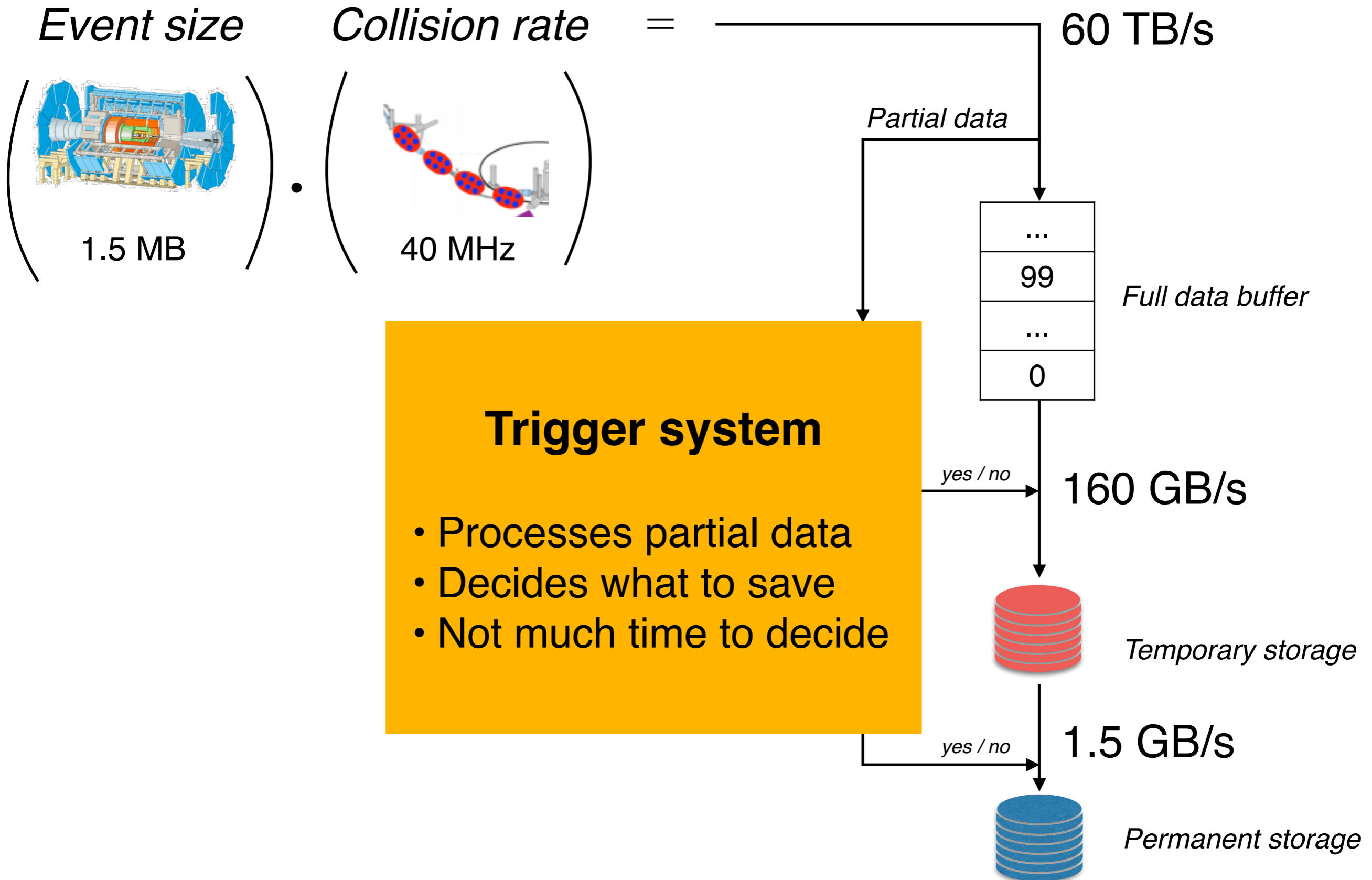  - Tree flattening & merging

- Firmware design
  - Bin Engines

- Physics results (simulated)
  - VBF Higgs vs. multijet

---

- Backup slides
  - Comparison to hls4ml
  - Test bench

*Event size*  •  *Collision rate*  =  ———————  60 TB/s

1.5 MB  •  40 MHz

*Partial data*

| ... |
|:---:|
| 99 |
| ... |
| 0 |

*Full data buffer*

## **Trigger system**

- Processes partial data
- Decides what to save
- Not much time to decide

*yes / no*  160 GB/s

*Temporary storage*

*yes / no*  1.5 GB/s

*Permanent storage*

**FW X Machina**    ***This talk***

65 TB/s

*Partial data*

| Level | Latency |
|-------|---------|

| ... |
|-----|
| 99 |
| ... |
| 0 |

**FPGA**

**Custom electronics**

**L1**    **O(1) µs**

Reduce rate to
100 kHz

*yes / no*    160 GB/s

**50k CPUs**

*Software*

HLT    O(1) s

Reduce rate to
1-10 kHz

*yes / no*    1.5 GB/s

Source: http://cern.ch/twiki/pub/Atlas/TDAQSpeakersCommitteeCommonReferences/tdaqFullNew2017.pdf

- **Workflow**

- Optimization

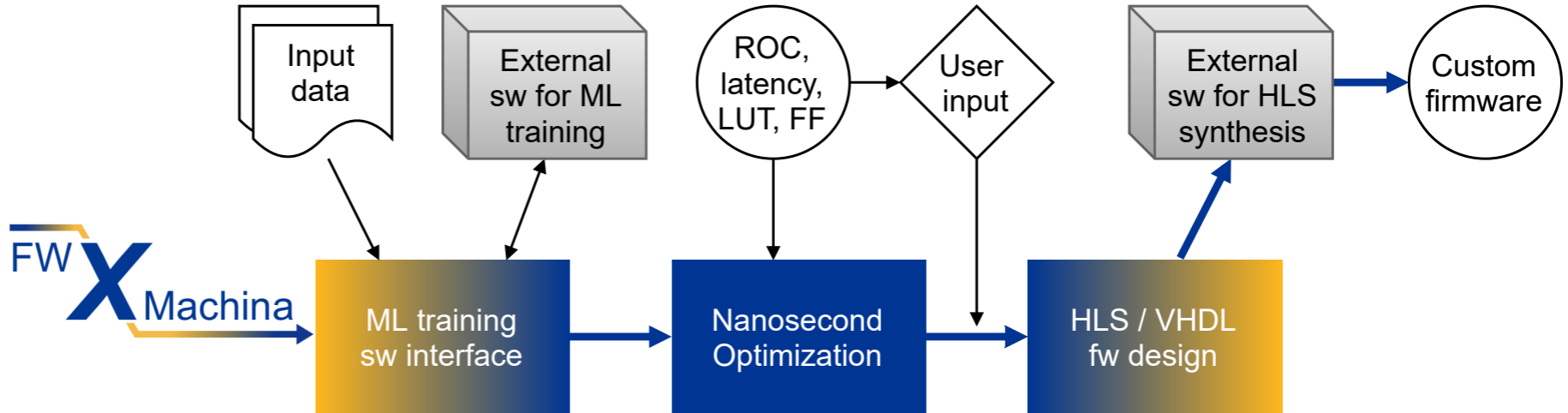- Use bit integer precision

- **Workflow**

- **Optimization**

- Use bit integer precision



- Will discuss next:

Tree Flattener

Forest Merger

- ## Workflow



- ## Optimization

- ## **Use bit integer precision**

    E.g., ap_int⟨8⟩ means the variable is represented by a range from 0 to 255.

    **Transformation**

    $$c_{\text{int}} = f(c_{\text{float}}) = \left\lfloor \frac{c_{\text{float}} - c_{\text{min}}}{c_{\text{max}} - c_{\text{min}}} \cdot \left(2^N - 1\right) \right\rfloor$$

    *Floor operation*

- ## Advantages & subtleties

    Bit integers represents a wide range without sacrificing float precision

    *Pre-evaluate f*     *Firmware only adds*

    $$f(x_1 + x_2) = f(x_1) + f(x_2)$$

    *Equal up to one bit because of floor*

**Conventional tree structure**

**2d plane: $x_a$ vs. $x_b$**

**Root node**

**First step**

**Depth i**

start

false    true

$q_i$: $x_a \geq c_i$
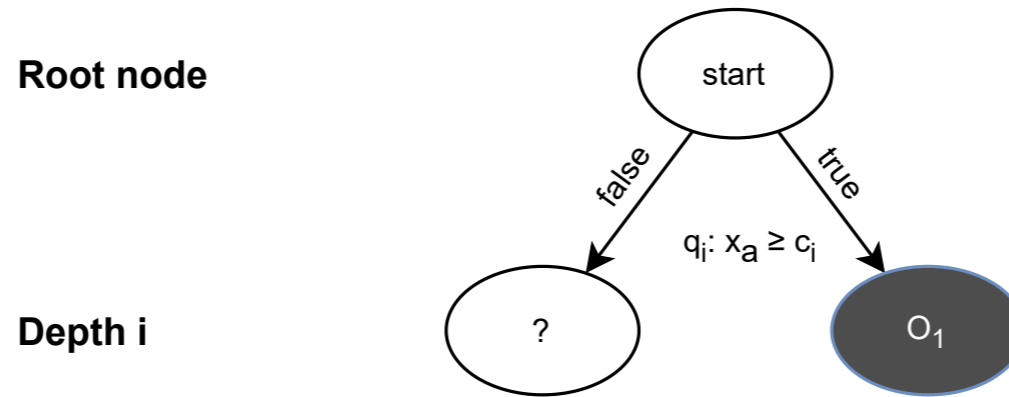
?    $O_1$

$x_b$

$c_{ii}$

$O_1$

$x_a$

$c_i$

- **Advantages** & subtleties
  - Cut thresholds & weights determined during training
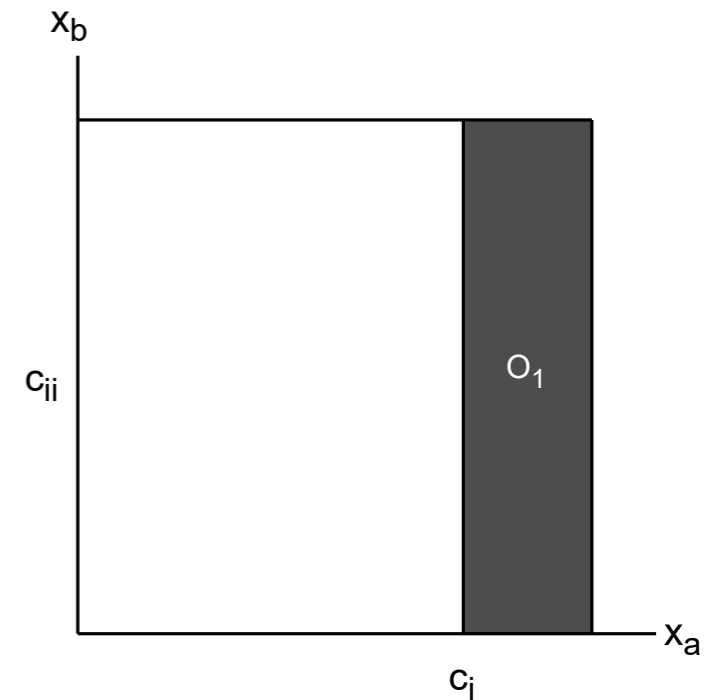  - Danger of "memorizing" boundaries (overtraining), so must consider a forest
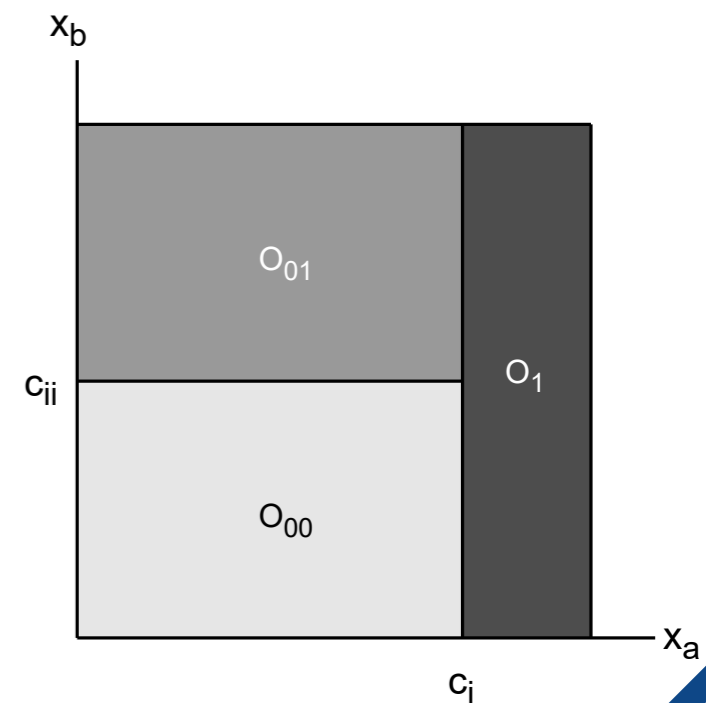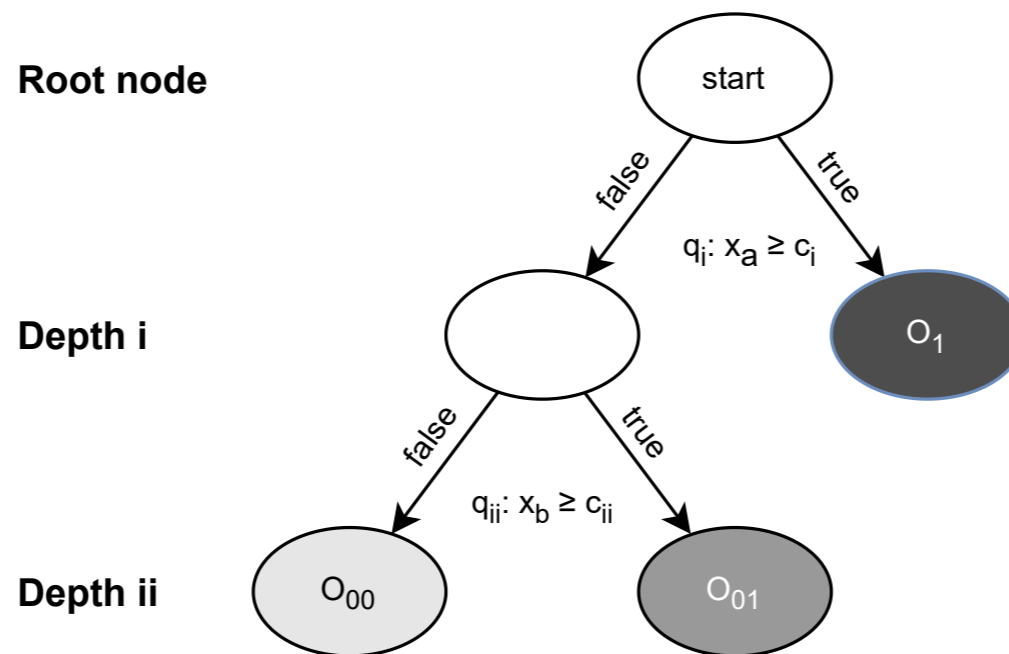
**Conventional tree structure**

**2d plane: $x_a$ vs. $x_b$**

**First step**

Root node

Depth i

start

false    true

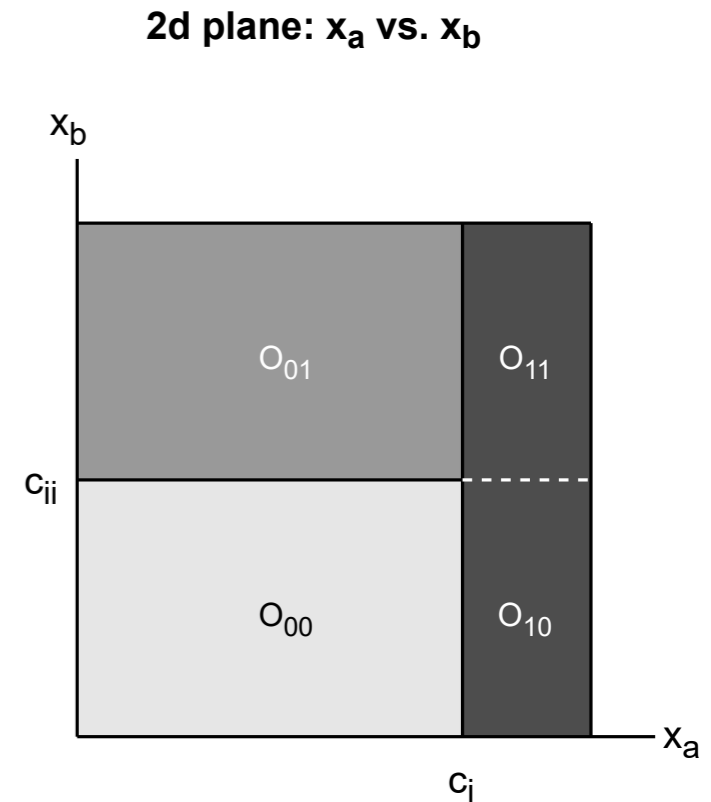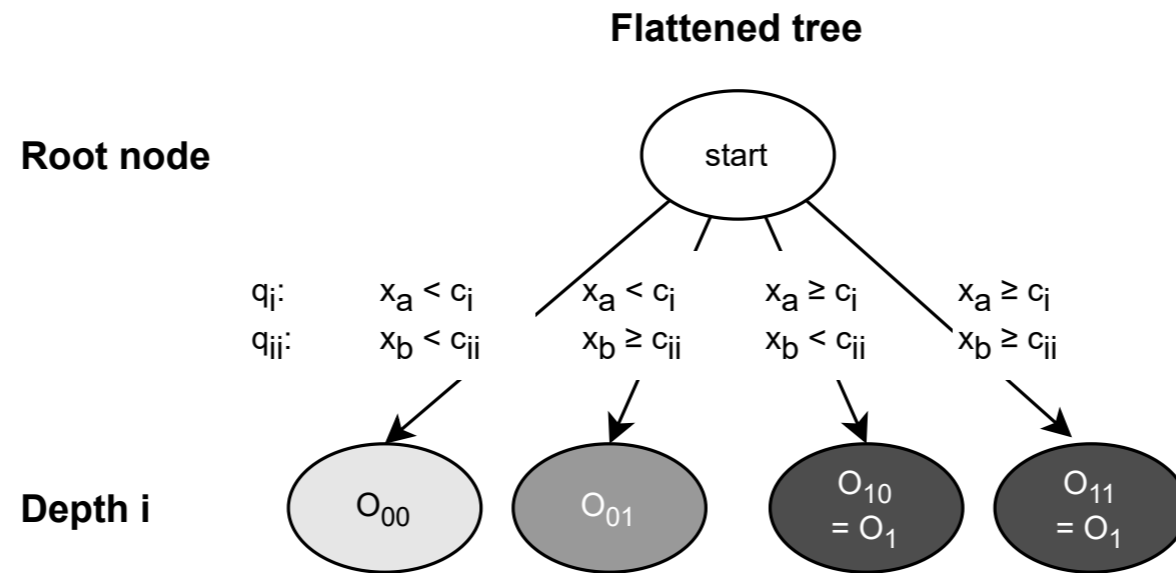$q_i$: $x_a \geq c_i$

?    $O_1$



- **Advantages & subtleties**
  - Deterministic, conventional style
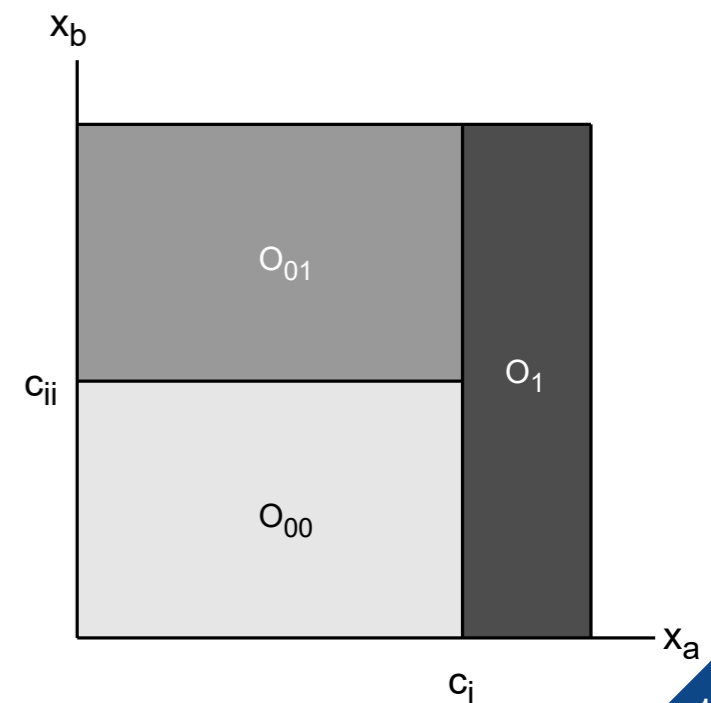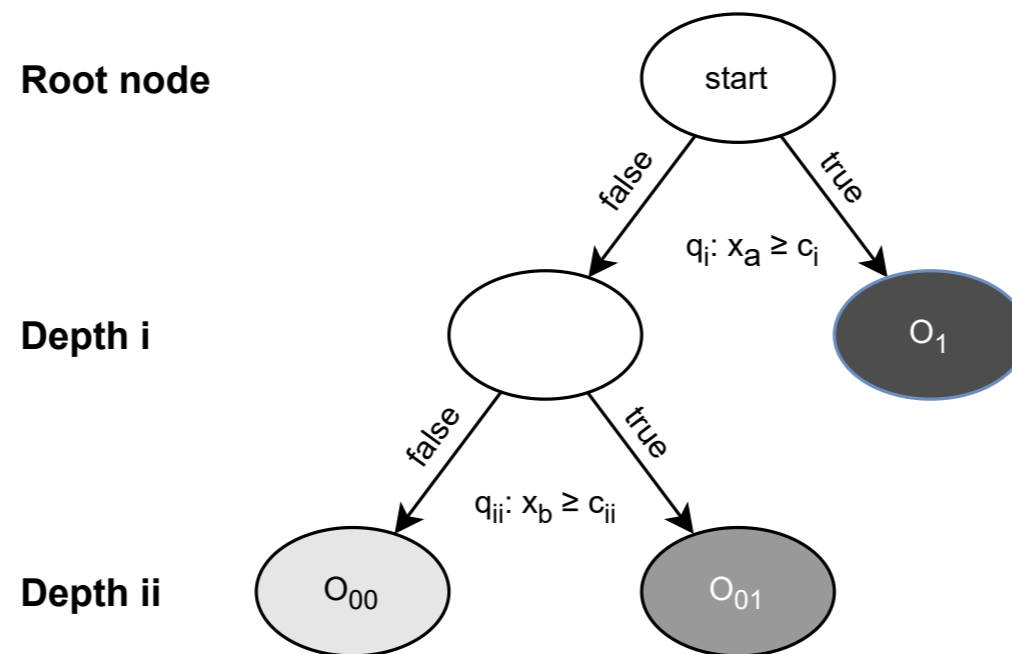  - Cuts in each axis is not independent of each other, so recursive

**Full tree**

Root node

Depth i

Depth ii

start

false    true

$q_i$: $x_a \geq c_i$

$O_1$

false    true

$q_{ii}$: $x_b \geq c_{ii}$

$O_{00}$    $O_{01}$

**Flattened tree**

**2d plane: $x_a$ vs. $x_b$**

**Root node**

Our approach

start

$q_i$:     $x_a < c_i$    $x_a < c_i$    $x_a \geq c_i$    $x_a \geq c_i$

$q_{ii}$:     $x_b < c_{ii}$    $x_b \geq c_{ii}$    $x_b < c_{ii}$    $x_b \geq c_{ii}$

**Depth i**

$O_{00}$     $O_{01}$     $O_{10} = O_1$     $O_{11} = O_1$

$x_b$

$O_{01}$    $O_{11}$

$c_{ii}$

$O_{00}$    $O_{10}$

$c_i$    $x_a$

- ## Advantages & subtleties
  - Each axis is independent of each other → Bin search problem on a grid
  - Does not scale well for very deep trees (but do you really need it at L1?)

**Root node**

Full tree

start

$q_i$: $x_a \geq c_i$

false    true

**Depth i**

$O_1$

false    true

$q_{ii}$: $x_b \geq c_{ii}$

**Depth ii**

$O_{00}$     $O_{01}$
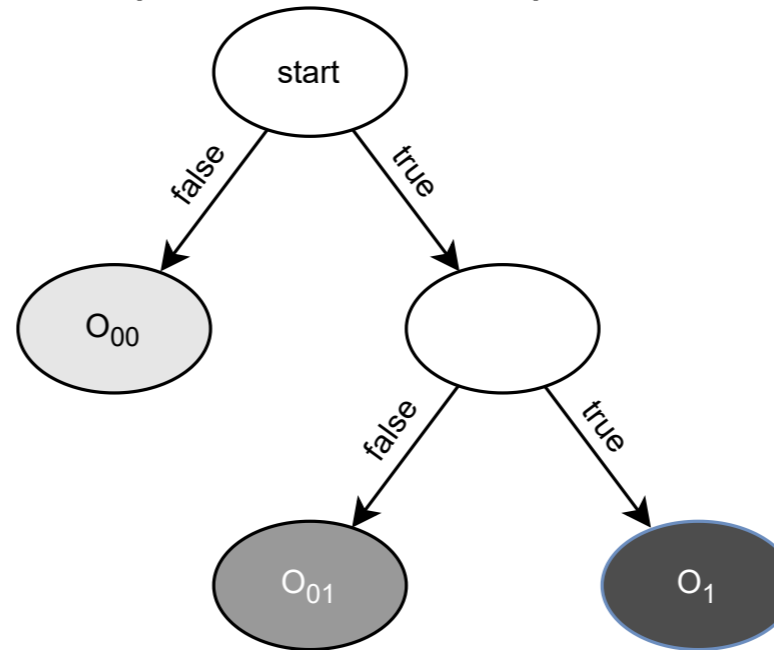
$x_b$

$O_{01}$

$c_{ii}$     $O_1$

$O_{00}$

$c_i$    $x_a$

- Advantages & subtleties
  - Use TMVA software to train the BDT (support for other sw coming)

Our approach

**2nd tree**

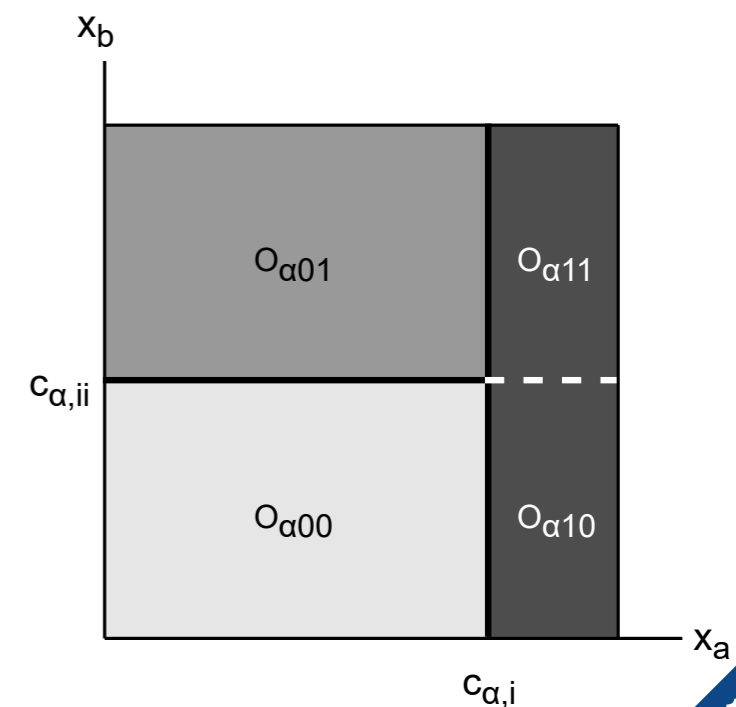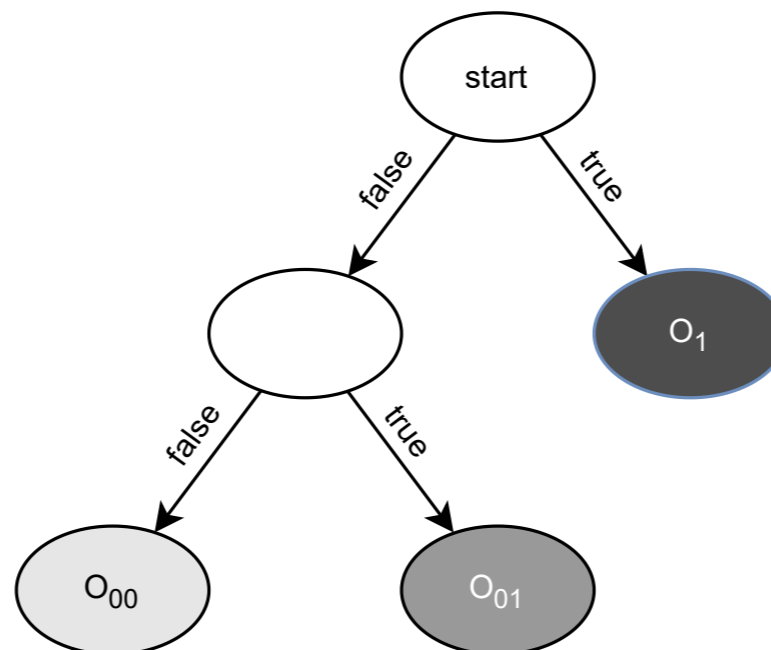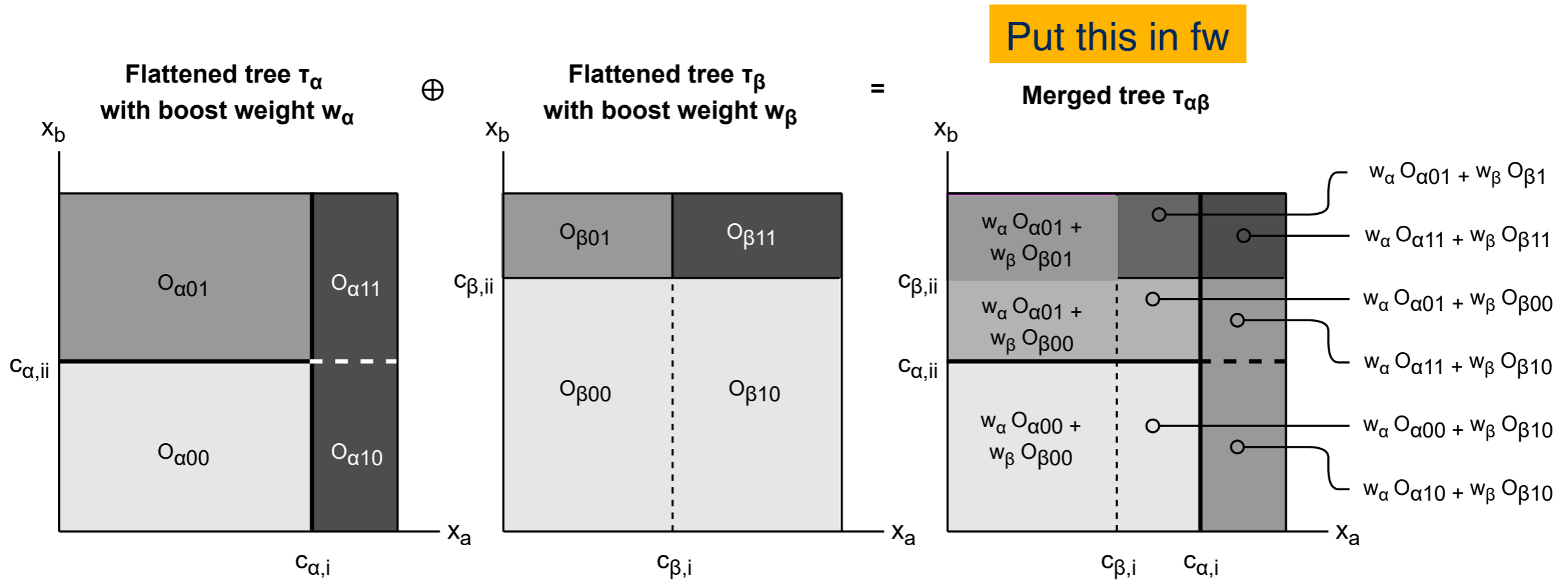**Decision tree $\tau_\beta$ with boost weight $w_\beta$**



- Can we pre-merge the trees for firmware? Yes, next slide.

**Decision tree $\tau_\alpha$ with boost weight $w_\alpha$**



**1st tree**

13

Put this in fw

Flattened tree $\tau_\alpha$ with boost weight $w_\alpha$ $\oplus$ Flattened tree $\tau_\beta$ with boost weight $w_\beta$ $=$ Merged tree $\tau_{\alpha\beta}$

- <span style="color:blue">Advantages</span> & <span style="color:red">subtleties</span>
  - <span style="color:blue">Merging is pre-processed before implementation in firmware</span>
  - <span style="color:red">This is using adaptive boosting. Gradient boosting cannot pre-merge, but we have approximations for that method to improve performance.</span>

- Physics impact of flattening & merging
  - None, bec. encodes the entirety of conventional approach
  - Firmware is a giant look-up table problem

TM Hong

- ## VBF Higgs vs. Multijet background
  - $\sigma_{Higgs}$ = 4 pb, two widely separated high-$p_T$ jets
  - $\sigma_{pp}$ = 80 mb, dominant process at LHC
  - Distributions given on the right

- ## We consider two decays of the Higgs
  - $H \to \nu\bar{\nu}\nu\bar{\nu}$, "invisible"
  - $H \to b\bar{b}b\bar{b}$, thru pseudoscalar decays

- ## Strategy
  - Train BDT to identify VBF jet pair,
    i.e., train BDT on Multijet vs. VBF $H \to \nu\bar{\nu}\nu\bar{\nu}$
  - Apply that BDT to Multijet vs. VBF $H \to b\bar{b}b\bar{b}$

- ## Why
  - If it works for VBF $H \to b\bar{b}b\bar{b}$, then it can be a trigger for VBF independent of the Higgs decay
  - Does it work? Next slide



15

- ## It works!
  - Reminder. Did *not* train on VBF H → b$\bar{b}$b$\bar{b}$
  - Subtlety re: jet selection (see paper)
  - Distributions given on the right

- ## Performance comparison
  - Try to mimic ATLAS HL-LHC cuts as best we can using Madgraph + Delphes
  - Two-fold signal efficiency improvement from ATLAS-inspired → fwX results

- ## Details
  - We validated our setup to reproduce the signal efficiency in the ATLAS Run-2 paper
  - Comparison using bit integers, not floats



Test samples
- Multijet
- VBF H → inv.
- VBF H → 4b

Notes
Trained vs.
VBF H → inv.
for all curves



VBF Higgs signal
- inv. BDT (float, all J)
- 4b BDT (float, all J)
- inv. BDT (bits, top 3)
- 4b BDT (bits, top 3)
- inv. cut (HL-LHC)
- 4b cut (HL-LHC)
- inv. cut (Run-2)
- 4b cut (Run-2)

BDT curves

3.2%  6.3%

Equal rate comparison of 4b

Better

- ## Ran two configurations
  - Optimized version
  - Non-optimized version (for comparison)
  - Both using 100 trees, max depth of 4
  - Results given on the right

- ## Performance
  - 5 clock ticks = 16 ns
  - Negligible resource usage

- ## Benchmark using e$^+$ vs. γ
  - In the paper, we also define <u>one set</u> of parameters to scale up <u>one param. at a time</u>
  - Uses 4 variables, 8 bits & same as above
  - 3 clock ticks = 10 ns
  - Negligible resource usage

|  | VBF H Optimized | VBF H Non-opt |
|---|---|---|
| $N_{var}$ | 5 | 7 |
| $N_{bit-var}$ | 8 | 12 |
| $N_{bit-score}$ | 16 | 16 |
| $N_{bin}$ | 40k | 1M |
| Latency | 5 ticks | 6 ticks |
| LUT | 1% | 1.5% |
| Flip Flops | ~0 | ~0 |
| BRAM | 2% | 30% |
| DSP | 0 | ~0 |

# Back up

- ## Details
  - Ideally we would run hls4ml's example & compare, but we can't as-is because they run a 5-class jet identification (b, W, top, g, q)
  - We ran hls4ml on the same dataset with the same configuration as in our paper

| Parameter | fwXmachina | hls4ml-Conifer | Comments |
|---|---|---|---|
| **ML training setup** | | | |
| Training software | *TMVA* | *TMVA* | same |
| Physics problem | electron vs. photon | electron vs. photon | same |
| Training samples | from ref. [56] | from ref. [56] | same |
| No. of event classes | 2 | 2 | same |
| No. of training trees | 100 | 100 | same |
| Max. depth | 4 | 4 | same |
| No. of input variables | 4 | 4 | See figure 18 |
| Other *TMVA* parameters | *TMVA* defaults | *TMVA* defaults | same |
| Nanosec. Optimization | Flattened & merged to 10 final trees, without TREE REMOVER or CUT ERASER | N/A | Unique to fwX |
| **FPGA and firmware setup** | | | |
| Chip family | Xilinx Virtex Ultrascale+ | Xilinx Virtex Ultrascale+ | same |
| Chip model | xcvu9p-flga2104-2L-e | xcvu9p-flga2104-2L-e | same |
| Vivado HLS version | 2019.2 | 2019.2 | same |
| Clock speed, period | 320 MHz, 3.125 ns | 320 MHz, 3.125 ns | same |
| Precision | ap_int⟨8⟩ | ap_ufixed⟨10, 5⟩ | See text |
| BIN ENGINE | BSBE | N/A | Unique to fwX |
| **FPGA cost** | | | |
| Latency | 3 clock ticks, 9.375 ns | 15 clock ticks, 46.875 ns | - |
| Interval | 1 clock tick, 3.125 ns | 1 clock tick, 3.125 ns | same |
| LUT | 1903, < 0.2% of total | 2.3 M, 192% of total | See caption |
| FF | 138, < 0.01% of total | 1.1 M, 44% of total | - |
| BRAM 18k | 8, < 0.2% of total | 0 | - |
| URAM | 0 | 0 | same |
| DSP | 0 | 0 | same |

- Same setup

- Comparison

| Parameter | Value | Comments |
|---|---|---|
| **FPGA setup** | | |
| Chip family | Xilinx Virtex Ultrascale+ | |
| Chip model | xcvu9p-flga2104-2L-e | |
| Vivado version | 2019.2.1 | |
| Synthesis type | C-Synthesis | |
| HLS or RTL | HLS | |
| HLS interface pragma | None | |
| Clock speed | 320 MHz | Clock period is 3.125 ns |
| **ML training configuration** | | |
| ML training method | Boosted decision tree | Binary classification |
| Boost method | Adaptive | AdaBoost with yes/no leaf |
| No. of event types to classify | 2 | Signal vs. background |
| No. of input variables | 4 | |
| No. of trees used for training | 100 | |
| Maximum tree depth | 4 | |
| **Nanosecond Optimization configuration** | | |
| Bin Engine type | Bit Shift Bin Engine (BSBE) | |
| No. of bits for input variables | 8 bits for each | |
| No. of bits for cut thresholds | 8 bits for each | |
| No. of bits for BDT output score | 8 bits | |
| No. of trees after merging | 10 | Tree Merger via ordered list |
| No. of final trees | 10, none removed | Tree Remover by truncation |
| No. of bins | 26132 | Cut Eraser not used |
| **FPGA cost** | | |
| Latency | 3 clock ticks | 9.375 ns |
| Interval | 1 clock tick | 3.125 ns |
| Look up tables | 1903 out of 1182240 | < 0.2% of available |
| Flip flops | 138 out of 2364480 | < 0.01% of available |
| Block RAM | 8 out of 4320 | < 0.2% of available |
| Ultra RAM | 0 out of 960 | - |
| Digital signal processors | 0 out of 6840 | - |



- 10 ns is independent of clock from 100-320 MHz

Table 9: List of input variables for the classification of the VBF Higgs boson vs. multijet process. Also given are the ATLAS-inspired cut-based offline thresholds for Run 2 [64] and HL-LHC [65]. For Run-2, differences arise with respect to the document when the $m_{jj}$ threshold is quoted as $1100\,\mathrm{GeV}$ for `L1 MJJ-500-NFF`; we use the $> 99\%$ offline efficiency point, which is achieved around $m_{jj} > 1300\,\mathrm{GeV}$. for others the offline thresholds are used. For HL-LHC, the single-level scheme values are quoted. The performance of the cut-based approach using these values is compared the performance to the BDT result in figure 16. The non-optimized (non-opt) configuration includes the five variables from the optimized configuration.

| Input variable | Description | ATLAS Run-2 offline cut [64], see caption | ATLAS HL-LHC offline cut [65], see caption | Used in BDT |
|---|---|---|---|---|
| $p_{\mathrm{T}1}$ | Leading jet $p_{\mathrm{T}}$ | $> 90\,\mathrm{GeV}$ | $> 75\,\mathrm{GeV}$ | - |
| $p_{\mathrm{T}2}$ | Subleading jet $p_{\mathrm{T}}$ | $> 80\,\mathrm{GeV}$ | $> 75\,\mathrm{GeV}$ | Optimized |
| $p_{\mathrm{T}12}$ | Sum $p_{\mathrm{T}1} + p_{\mathrm{T}2}$ | - | - | Optimized |
| $|\eta_1|$ | Leading jet $\eta$ | $< 3.2$ | - | - |
| $|\eta_2|$ | Subleading jet $\eta$ | $< 4.9$ | - | - |
| $\prod_\eta$ | Product $\eta_1 \cdot \eta_2$ | - | - | Optimized |
| $|\Delta\eta|$ | Separation in $|\eta_2 - \eta_1|$ | $> 4.0$ | $> 2.5$ | - |
| $|\Delta\phi|$ | Separation in $|\phi_2 - \phi_1|$ | $< 2.0$ | $< 2.5$ | non-opt |
| $|\Delta R|$ | $\sqrt{(\Delta\eta)^2 + (\Delta\phi)^2}$ | - | - | non-opt |
| $m_{jj}$ | Dijet invariant mass | $> 1300\,\mathrm{GeV}$ | - | Optimized |
| $p_T^{jj}$ | Dijet $p_{\mathrm{T}}$ | - | - | Optimized |

- Each variable is processed independently of each other

- Look up thresholds in memory, compare

- Bit shift to localize data
  - This is *fast*

- Use combinatoric logic as much as possible without multiplication. No explicit clocked operations.

- No difference seen wrt software implementation

# More info

All info at
https://fwx.pitt.edu



data set

git repo

doxygen

**fwXmachina Project**

## Welcome!

Information regarding the **fwX** project will be available on this page. This project is developed by members of the Hong Group in the Department of Physics and Astronomy.

## Where to find information

- The pre-print is available on the arXiv [2104.03408]
- Data samples used for the VBF Higgs and multijet study
    - Stephen Roche, Benjamin Carlson, Tae Min Hong (2021), *fwXmachina example: VBF Higgs vs multijet*, Mendeley Data, V1, doi: 10.17632/kp3myh3v89.1
- Download code
    - v1.0.0: https://gitlab.com/PittHongGroup/fwX/-/tree/v1.0.0 and Doxygen documentation
- COMING SOON
    - User Guide
    - Tutorials

## Communicate with us

- Sign-up for the announcements at Mailman to receive emails from fwx@list.pitt.edu.
- Send inquiries to fwx-developers@list.pitt.edu.

examples

code

driver file

**install**