



# Graph Neural Network(GNN) Inference on FPGA

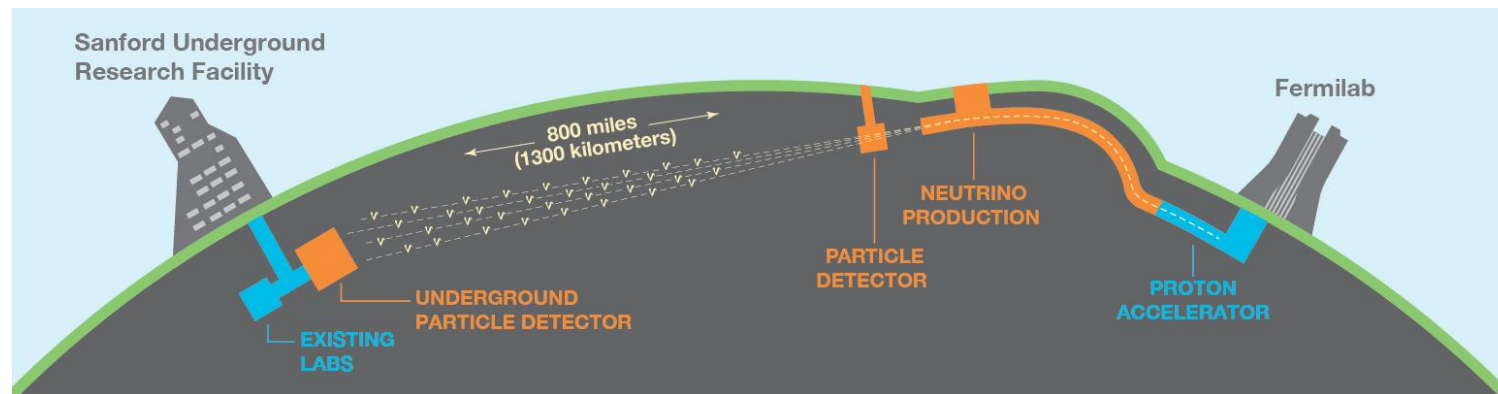
*CERN openlab Lightning Talks*

Kazi Ahmed Asif Fuad

Supervisor: **Sofia Vallecorsa**

15/08/2019

# Project Background



# Our Objective

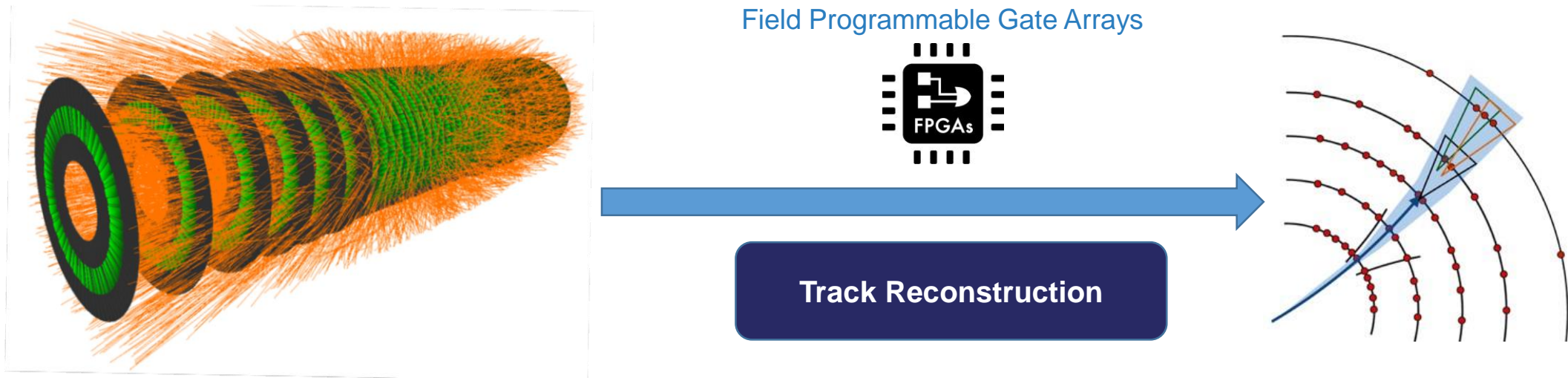


Image Based Methods  Space-Point Representation

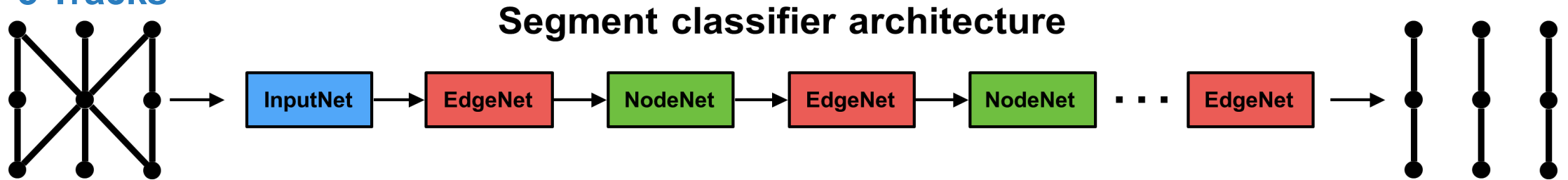
<https://indico.cern.ch/event/753577/contributions/3123602/attachments/1707996/2752966/acts-gnn-Aug30.pdf>

[https://indico.cern.ch/event/658267/contributions/2881175/attachments/1621912/2581064/Farell\\_heptrkx\\_ctd2018.pdf](https://indico.cern.ch/event/658267/contributions/2881175/attachments/1621912/2581064/Farell_heptrkx_ctd2018.pdf)

HEP.TrkX: <https://heptrkx.github.io/>

# Graph Neural Network (GNN)

3 Layers  
3 Tracks



With each iteration, the model propagates information through the graph, strengthens important connections, and weakens useless ones.

InputNet: New Features  
tanh activations

1 Layer MLP

EdgeNet: Edge Weights  
tanh activations  
sigmoid activation

2 Layer MLP

NodeNet: New Features  
tanh activations  
tanh activations

2 Layer MLP

<https://indico.cern.ch/event/753577/contributions/3123602/attachments/1707996/2752966/acts-gnn-Aug30.pdf>

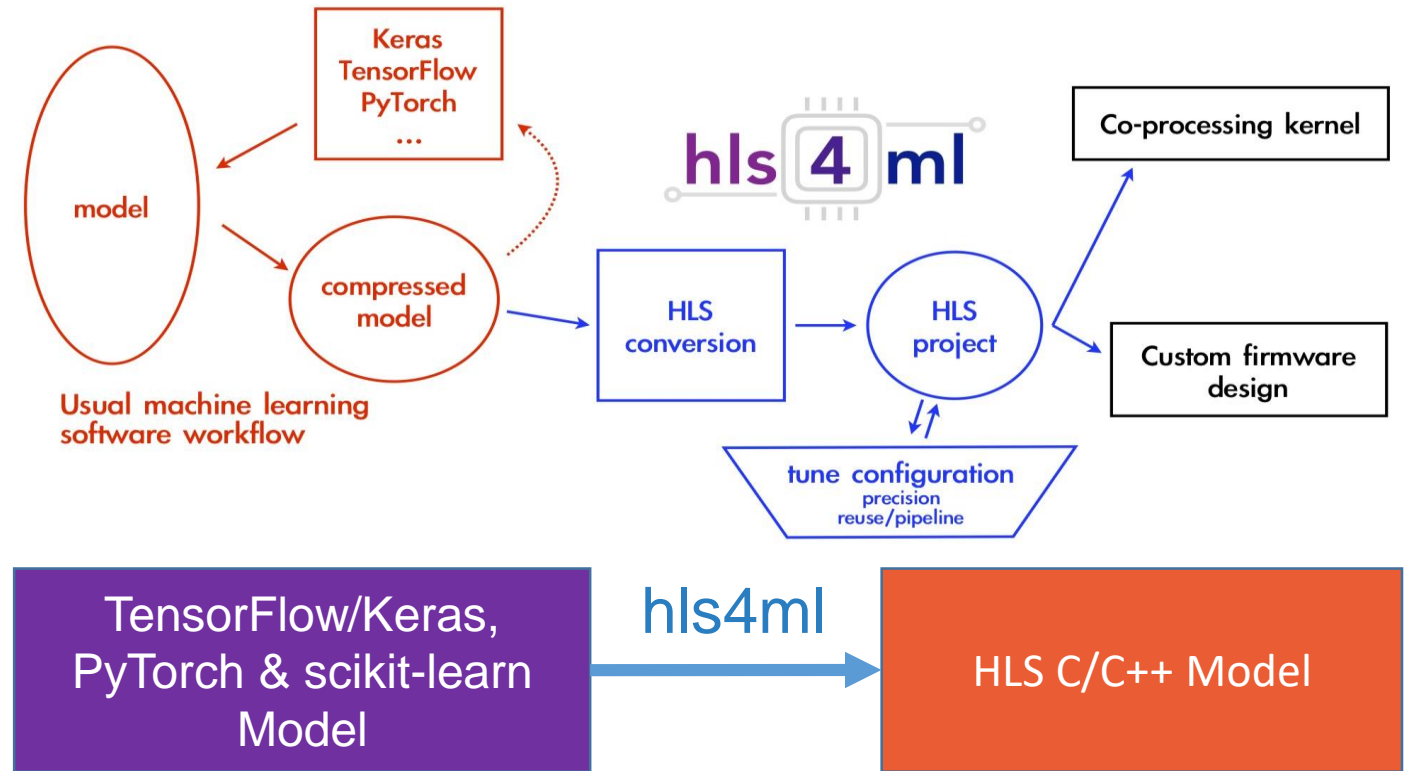
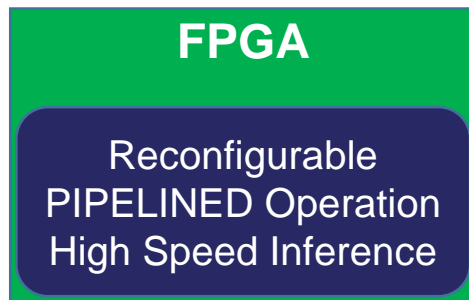
[https://indico.cern.ch/event/658267/contributions/2881175/attachments/1621912/2581064/Farrell\\_heptrkx\\_ctd2018.pdf](https://indico.cern.ch/event/658267/contributions/2881175/attachments/1621912/2581064/Farrell_heptrkx_ctd2018.pdf)

# Implementation on FPGA

High Level Synthesis



A package for machine learning inference in FPGAs.



hls4ml: <https://hls-fpga-machine-learning.github.io/hls4ml/>

# More FPGA Facts

- ✓ Basic Building Resource Blocks: LUTS, DSPs, Flip-Flops & BRAMs.
- ✓ Resource Utilization needs to be less than 100% to **fit** a design into FPGA.
- ✓ Resource Utilization of SLR less than 100% is good for the design.
- ✓ Reuse Factor means how many times the DSP(Multiplier + Adder) block will be used.
- ✓ PIPELINE Architecture is faster than Dataflow Architecture but utilizes more resources.

# My HLS Implementation

For HLS implementation, I have merged following implementations.

GNN implementation of **Javier M. G. Duarte (Fermilab)**

<https://github.com/hls-fpga-machine-learning/hls4ml/tree/jmgd/graph/example-prjs/graph>

Reference for GNN

+

Large Dense Layers Implementation from **Vladimir Loncar (CERN)**:

<https://github.com/vloncar/hls4ml/tree/hack6>

Reference for NN

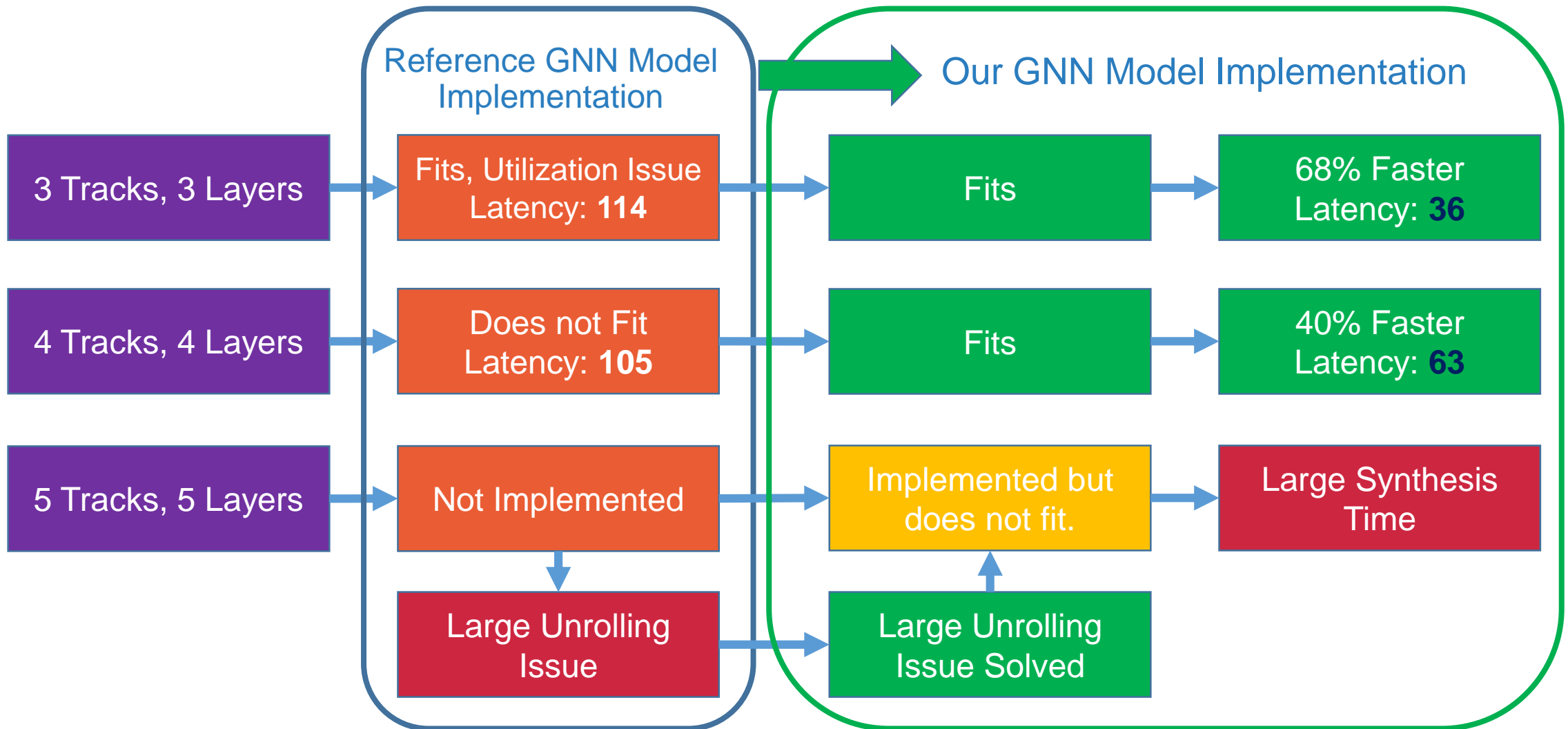


My implementations are available at:

<https://github.com/belloworld/hls4ml/tree/hack6/example-prjs/GNN>

Our Implementation

# Results for Pipeline Architecture





# Issues, We are Facing in Pipeline

- ❑ Reuse Factor Not Working
- ❑ Large Synthesis Time

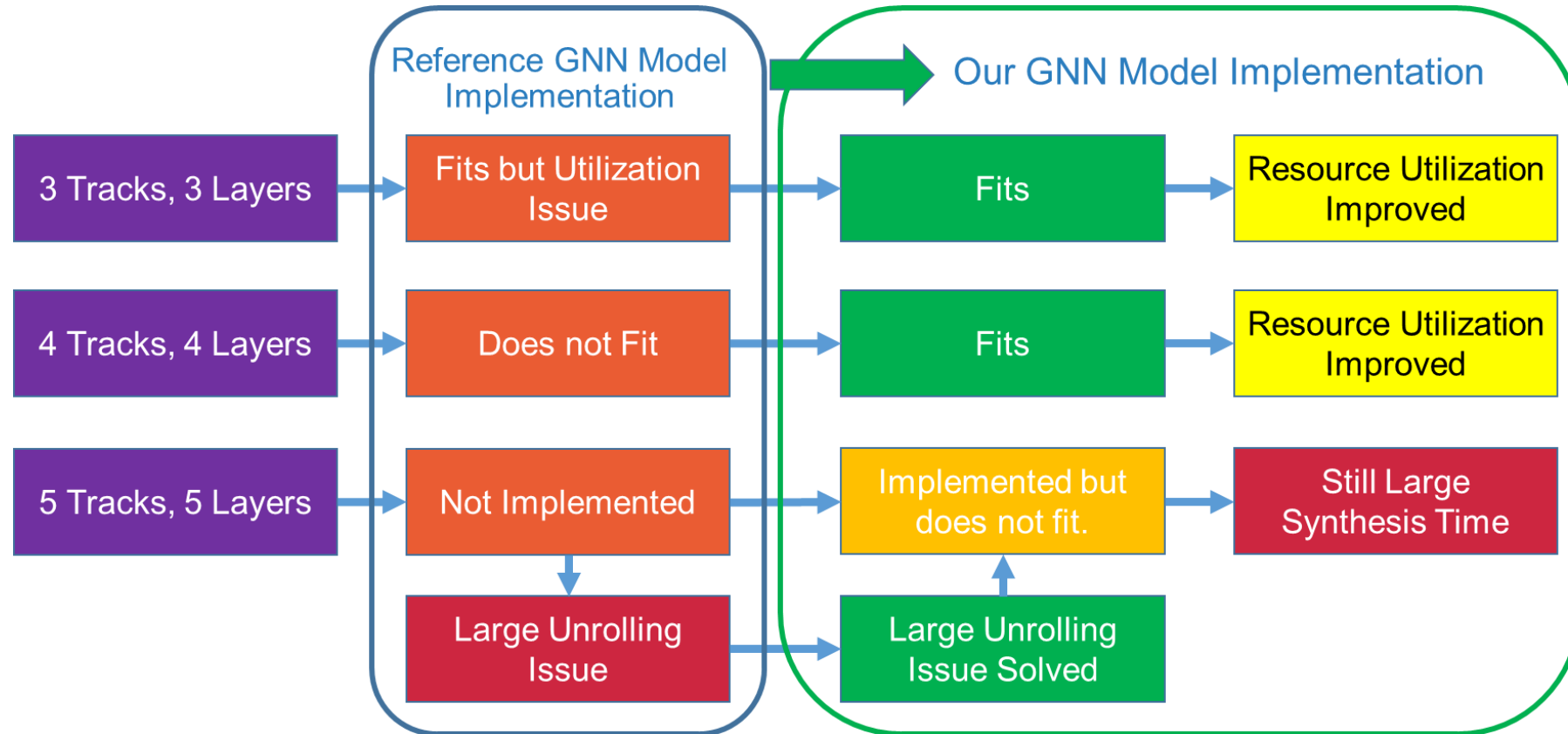
After discussions.....

Opting to.....

DATAFLOW

# Results for Dataflow Architecture

REUSE Factor Works but **Long Synthesis time not solved yet!**



# Things to do and Future Work(!)

In summary,

My 1<sup>st</sup> implemented GNNs around around **40% faster** in Pipeline architecture.

My 2<sup>nd</sup> implementations are using around **45% less resources** than the reference.

- More Investigation on the design.
- Different perspective for large unrolling issue
- Run the 3 Tracks, 3 Layers GNN on Kintex FPGA.

Ultimate Target is the 10 Tracks, 10 Layers GNN (!)

# Special Thanks to...



Sofia  
Vallecorsa



Vladimir  
Loncar





# QUESTIONS?

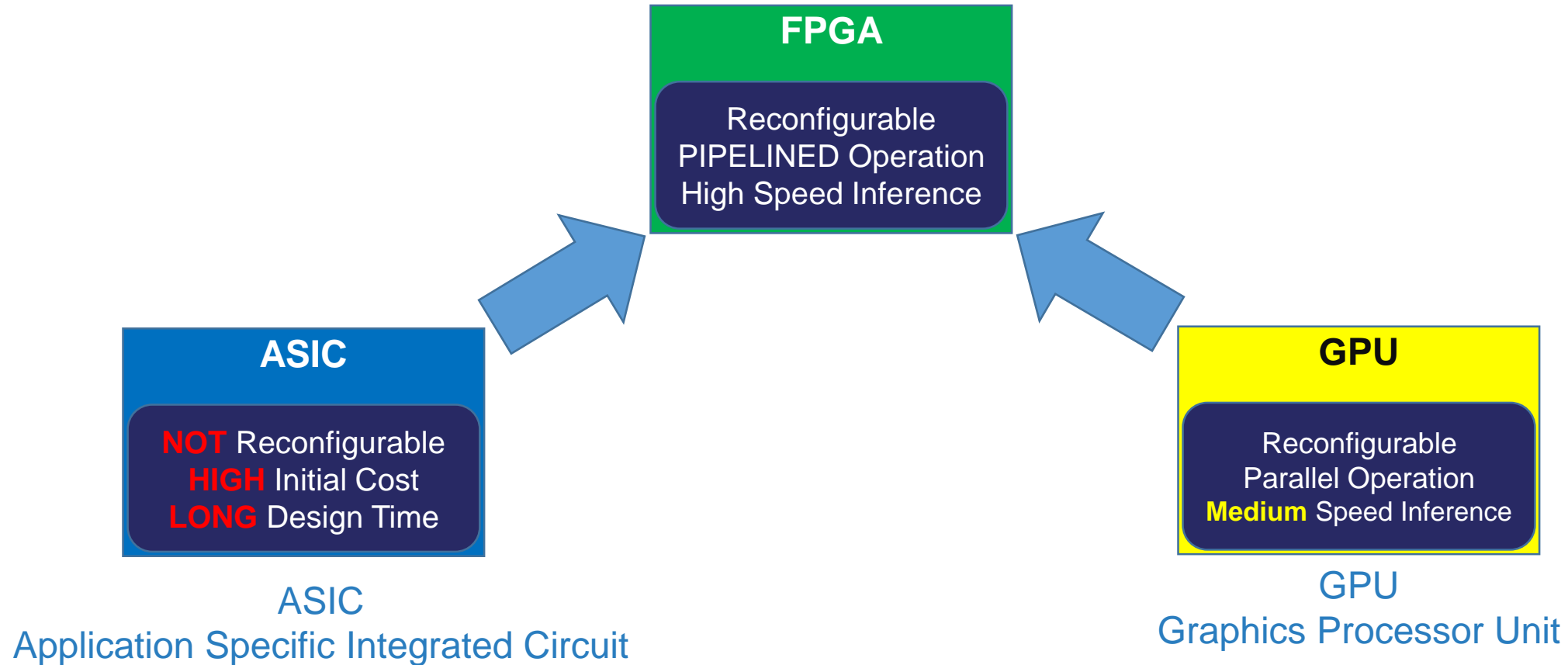
 [asif.ahmed.fuad@gmail.com](mailto:asif.ahmed.fuad@gmail.com)

 <https://www.linkedin.com/in/asif-fuad/>

# Additional Slides

# Why FPGA?

FPGA  
Field Programmable Gate Arrays



<https://www.arrow.com/en/research-and-events/articles/fpga-vs-cpu-vs-gpu-vs-microcontroller>

<https://lancesimms.com/Microprocessors/CPU vs GPU vs FPGA.html>

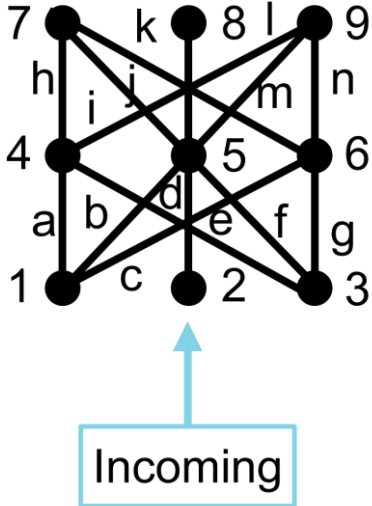
<https://numato.com/blog/differences-between-fpga-and-asics/>

# A Simple Graph

A Simple 3 Layers Graph

Each layer has 3 Nodes(hits)

Our objective is to identify “Good” segments

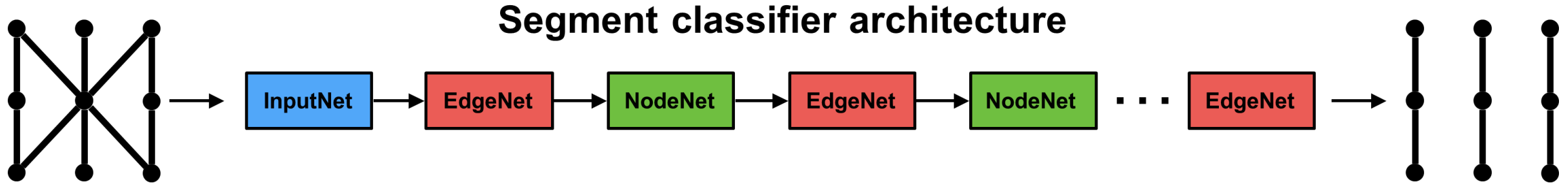


<https://indico.cern.ch/event/753577/contributions/3123602/attachments/1707996/2752966/acts-gnn-Aug30.pdf>

[https://indico.cern.ch/event/658267/contributions/2881175/attachments/1621912/2581064/Farrell\\_heptrkx\\_ctd2018.pdf](https://indico.cern.ch/event/658267/contributions/2881175/attachments/1621912/2581064/Farrell_heptrkx_ctd2018.pdf)



# Graph Neural Network (GNN)



*With each iteration, the model propagates information through the graph, strengthens important connections, and weakens useless ones.*

## Definitions

$X$   $\longrightarrow$  (N x D) node feature matrix

$R_i$   $\longrightarrow$  (N x E) association matrix of nodes to input edges

$R_o$   $\longrightarrow$  (N x E) association matrix of nodes to output edges

<https://indico.cern.ch/event/753577/contributions/3123602/attachments/1707996/2752966/acts-gnn-Aug30.pdf>

[https://indico.cern.ch/event/658267/contributions/2881175/attachments/1621912/2581064/Farrell\\_heptrkx\\_ctd2018.pdf](https://indico.cern.ch/event/658267/contributions/2881175/attachments/1621912/2581064/Farrell_heptrkx_ctd2018.pdf)

HEP.TrkX: <https://heptrkx.github.io/>

# Graph Neural Network (GNN)

- The edge network is a 2-layer MLP with tanh and sigmoid activations:

$$w = f_{\text{edge}}(R_i^T X, R_o^T X) \longrightarrow \text{(E) edge weight array}$$

- The node network is a 2-layer MLP with tanh activations:

$$X' = f_{\text{node}}\left((R_i \odot w)R_o^T X, (R_o \odot w)R_i^T X, X\right) \longrightarrow \text{(N x D) node features}$$

- Outputs
  - Node classifier: binary classifier layer gives score for each node
  - Segment classifier: final edge network application gives score for each edge

<https://indico.cern.ch/event/753577/contributions/3123602/attachments/1707996/2752966/acts-gnn-Aug30.pdf>

[https://indico.cern.ch/event/658267/contributions/2881175/attachments/1621912/2581064/Farrell\\_heptrkx\\_ctd2018.pdf](https://indico.cern.ch/event/658267/contributions/2881175/attachments/1621912/2581064/Farrell_heptrkx_ctd2018.pdf)

HEP.TrkX: <https://heptrkx.github.io/>

# 4 Tracks, 4 Layers 1 Iteration

## Input Network:

12 weights

## Edge Network:

60 weights

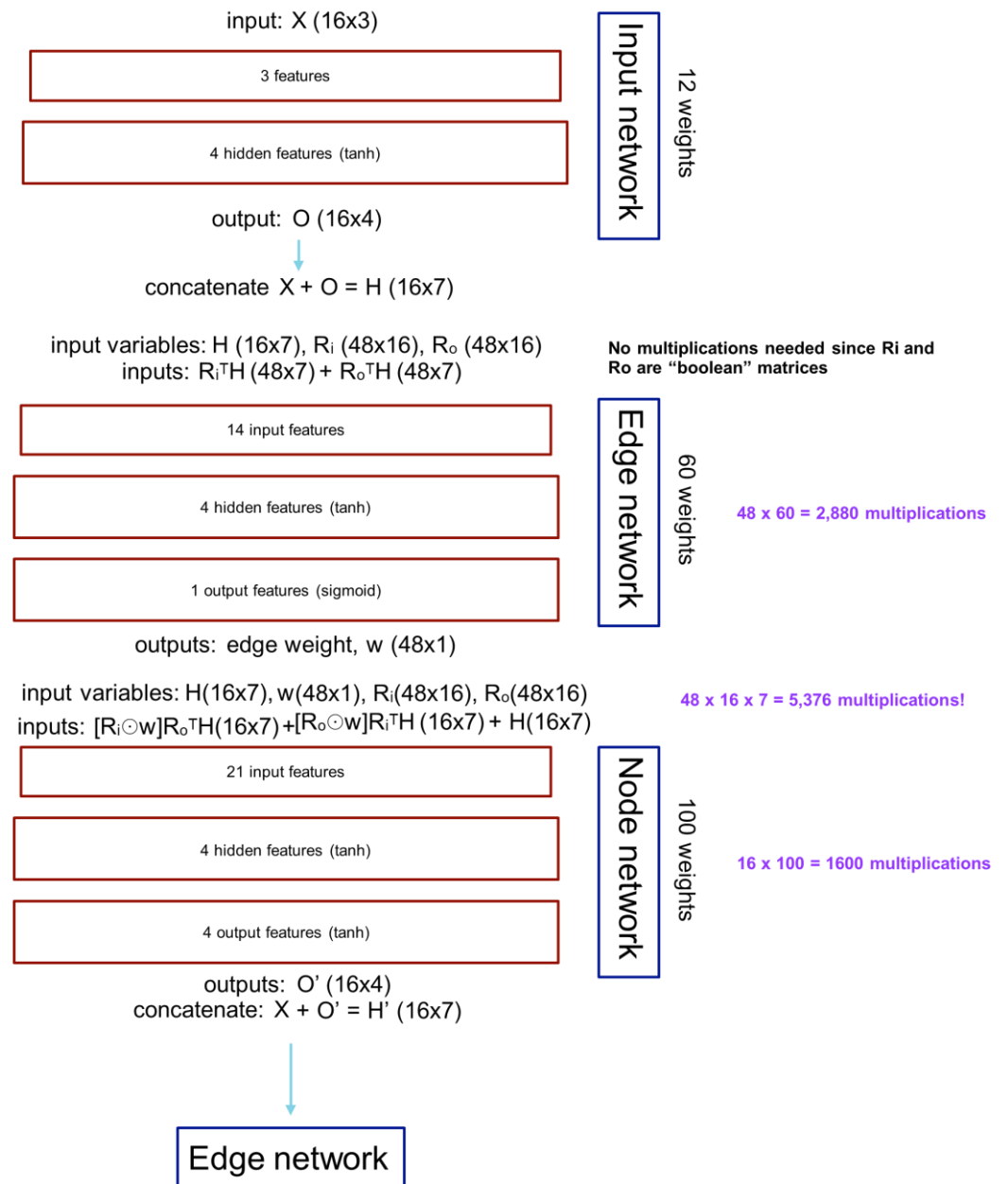
$48 \times 60 = 2880$  multiplications

$48 \times 16 \times 7 = 5,376$  multiplications

## Node Network:

100 weights

$16 \times 100 = 1600$  multiplications



# TIMING & RESOURCE USAGE (3 TRACKS, 3 LAYERS)

*Kintex FPGA: xcku115-flva1517-1-c : Pipeline Architecture*

GNN Resources Usage for Pipeline Architecture		Device: xcku115-flva1517-1-c							
		Utilization Estimates: Vivado HLS C Synthesis				Utilization Estimates: Vivado Synthesis			
		DSP48E	Change	LUT	Change	DSP48E	Change	CLB LUT	Change
	Available	<b>5520</b>	<b>na</b>	<b>663360</b>	<b>na</b>	<b>5520</b>	<b>na</b>	<b>663360</b>	<b>na</b>
	Available SLR	<b>2760</b>	<b>na</b>	<b>331680</b>	<b>na</b>	<b>na</b>	<b>na</b>	<b>na</b>	<b>na</b>
Reuse=1	Total(Used)	<b>5067</b>	<b>-776.64%</b>	<b>420412</b>	<b>-15.90%</b>	<b>5049</b>	<b>-773.53%</b>	<b>143112</b>	<b>60.55%</b>
	Utilization(%)	<b>91</b>	<b>-810.00%</b>	<b>63</b>	<b>-16.67%</b>	<b>91.47</b>	<b>-814.70%</b>	<b>21.57</b>	<b>60.06%</b>
	Utilization SLR (%)	<b>183</b>	<b>-815.00%</b>	<b>126</b>	<b>-15.60%</b>	<b>na</b>		<b>na</b>	
Reuse=7	Total(Used)	<b>1484</b>	<b>-156.75%</b>	<b>295398</b>	<b>18.56%</b>	<b>3309</b>	<b>-472.49%</b>	<b>106199</b>	<b>70.72%</b>
	Utilization(%)	<b>26</b>	<b>-160.00%</b>	<b>44</b>	<b>18.52%</b>	<b>59.95</b>	<b>-499.50%</b>	<b>16.01</b>	<b>70.35%</b>
	Utilization SLR (%)	<b>53</b>	<b>-165.00%</b>	<b>89</b>	<b>18.35%</b>				
Reuse=21	Total(Used)	<b>1161</b>	<b>-100.87%</b>	<b>285845</b>	<b>21.19%</b>	<b>3023</b>	<b>-423.01%</b>	<b>107392</b>	<b>70.39%</b>
	Utilization(%)	<b>21</b>	<b>-110.00%</b>	<b>43</b>	<b>20.37%</b>	<b>54.76</b>	<b>-447.60%</b>	<b>16.19</b>	<b>70.02%</b>
	Utilization SLR (%)	<b>42</b>	<b>-110.00%</b>	<b>86</b>	<b>21.10%</b>				

	Latency (Clock Cycles)			
	Latency			
	Min	Change	Max	Change
Reuse=1	<b>21</b>	<b>81.58%</b>	<b>21</b>	<b>81.58%</b>
Reuse=7	<b>36</b>	<b>68.42%</b>	<b>36</b>	<b>68.42%</b>
Reuse=21	<b>73</b>	<b>35.96%</b>	<b>73</b>	<b>35.96%</b>

# TIMING & RESOURCE USAGE (4 TRACKS, 4 LAYERS)

*Kintex FPGA: xcku115-flva1517-1-c : Pipeline Architecture*

GNN Resources Usage for Pipeline Architecture		Device: xcku115-flva1517-1-c							
		Utilization Estimates: Vivado HLS C Synthesis				Utilization Estimates: Vivado Synthesis			
		DSP48E	Change	LUT	Change	DSP48E	Change	CLB LUT	Change
	Available	<b>5520</b>	<b>na</b>	<b>663360</b>	<b>na</b>	<b>5520</b>	<b>na</b>	<b>663360</b>	<b>na</b>
	Available SLR	<b>2760</b>	<b>na</b>	<b>331680</b>	<b>na</b>	<b>na</b>	<b>na</b>	<b>na</b>	<b>na</b>
Reuse=1	Total(Used)	<b>17616</b>	<b>-823.27%</b>	<b>1798687</b>	<b>-19.40%</b>	<b>5520</b>	<b>-189.31%</b>	<b>2285042</b>	<b>-51.68%</b>
	Utilization(%)	<b>319</b>	<b>-838.24%</b>	<b>271</b>	<b>-19.38%</b>	<b>100</b>	<b>-194.12%</b>	<b>344.46</b>	<b>-51.74%</b>
	Utilization SLR (%)	<b>638</b>	<b>-824.64%</b>	<b>542</b>	<b>-19.38%</b>				
Reuse=7	Total(Used)	<b>5664</b>	<b>-196.86%</b>	<b>1386769</b>	<b>7.95%</b>	<b>2432</b>	<b>-27.46%</b>	<b>1181929</b>	<b>21.54%</b>
	Utilization(%)	<b>102</b>	<b>-200.00%</b>	<b>209</b>	<b>7.93%</b>	<b>44.06</b>	<b>-29.59%</b>	<b>178.17</b>	<b>21.51%</b>
	Utilization SLR (%)	<b>205</b>	<b>-197.10%</b>	<b>418</b>	<b>7.93%</b>				
Reuse=21	Total(Used)	<b>4640</b>	<b>-143.19%</b>	<b>1355382</b>	<b>10.03%</b>	<b>2582</b>	<b>-35.32%</b>	<b>1025563</b>	<b>31.92%</b>
	Utilization(%)	<b>84</b>	<b>-147.06%</b>	<b>204</b>	<b>10.13%</b>	<b>46.78</b>	<b>-37.59%</b>	<b>154.6</b>	<b>31.89%</b>
	Utilization SLR (%)	<b>168</b>	<b>-143.48%</b>	<b>408</b>	<b>10.13%</b>				

	Latency (Clock Cycles)			
	Latency			
	Min	Change	Max	Change
Reuse=1	<b>23</b>	<b>78.10%</b>	<b>23</b>	<b>78.10%</b>
Reuse=7	<b>32</b>	<b>69.52%</b>	<b>32</b>	<b>69.52%</b>
Reuse=21	<b>63</b>	<b>40.00%</b>	<b>63</b>	<b>40.00%</b>

# TIMING & RESOURCE USAGE (4 TRACKS, 4 LAYERS)

*Virtex FPGA: xcvu13p-fhga2104-1-i : Pipeline Architecture*

GNN Resources Usage for Pipeline Architecture		Device: xcvu13p-fhga2104-1-i			
		Vivado HLS C Synthesis		Vivado Synthesis	
		DSP48E	LUT	DSP48E	CLB LUT
	Available	12288	1728000	12288	1728000
	Available SLR	na	na	na	na
Reuse=1	Total(Used)	<b>17616</b>	<b>1790783</b>	<b>12288</b>	<b>991030</b>
	Utilization(%)	<b>143</b>	<b>103</b>	<b>100</b>	<b>57.35</b>
	Utilization SLR (%)	na	na	na	na
Reuse=7	Total(Used)	<b>5664</b>	<b>1380016</b>	<b>12272</b>	<b>654134</b>
	Utilization(%)	<b>46</b>	<b>79</b>	<b>99.87</b>	<b>37.85</b>
	Utilization SLR (%)	na	na	na	na
Reuse=21	Total(Used)	<b>4640</b>	<b>1355242</b>	<b>12272</b>	<b>629730</b>
	Utilization(%)	<b>37</b>	<b>78</b>	<b>99.87</b>	<b>36.44</b>
	Utilization SLR (%)	na	na	na	na

	Latency (Clock Cycles)	
	Latency	
	Min	Max
Reuse=1	<b>21</b>	<b>21</b>
Reuse=7	<b>29</b>	<b>29</b>
Reuse=21	<b>61</b>	<b>61</b>