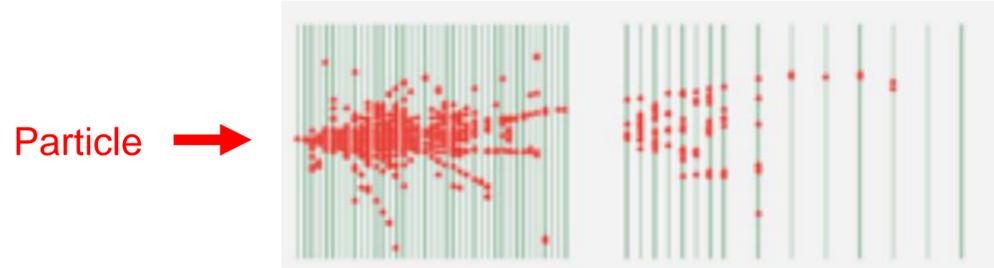


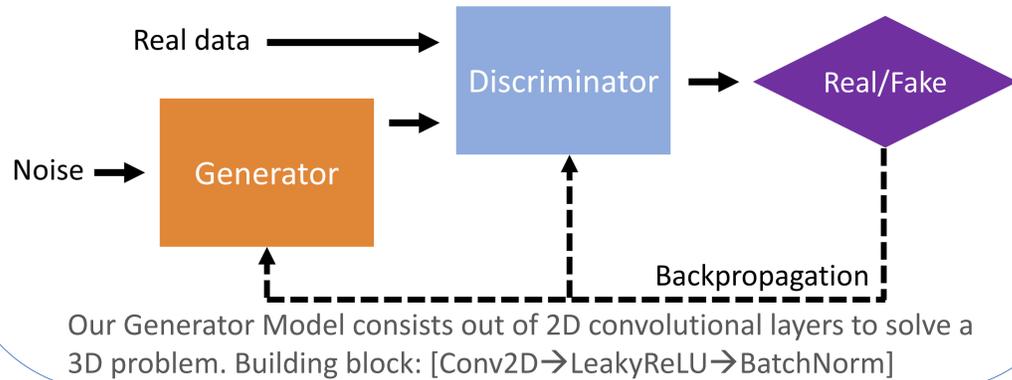
# Reduced Precision Strategies for Deep Learning: A GAN Use Case from High Energy Physics

## Electromagnetic Calorimeter Simulations



## Generative Adversarial Networks (GANs)

to replace traditional Monte Carlo Geant4 simulations



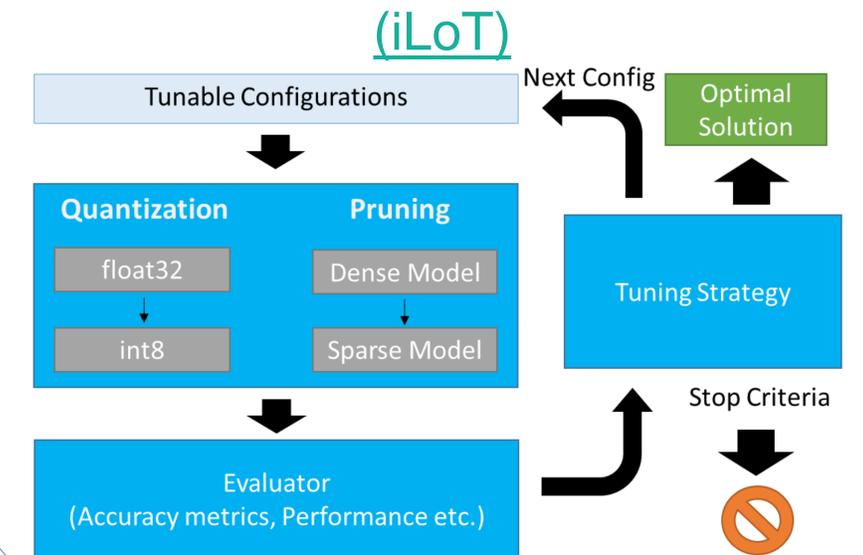
## Quantization

Converting a number from higher to lower precision

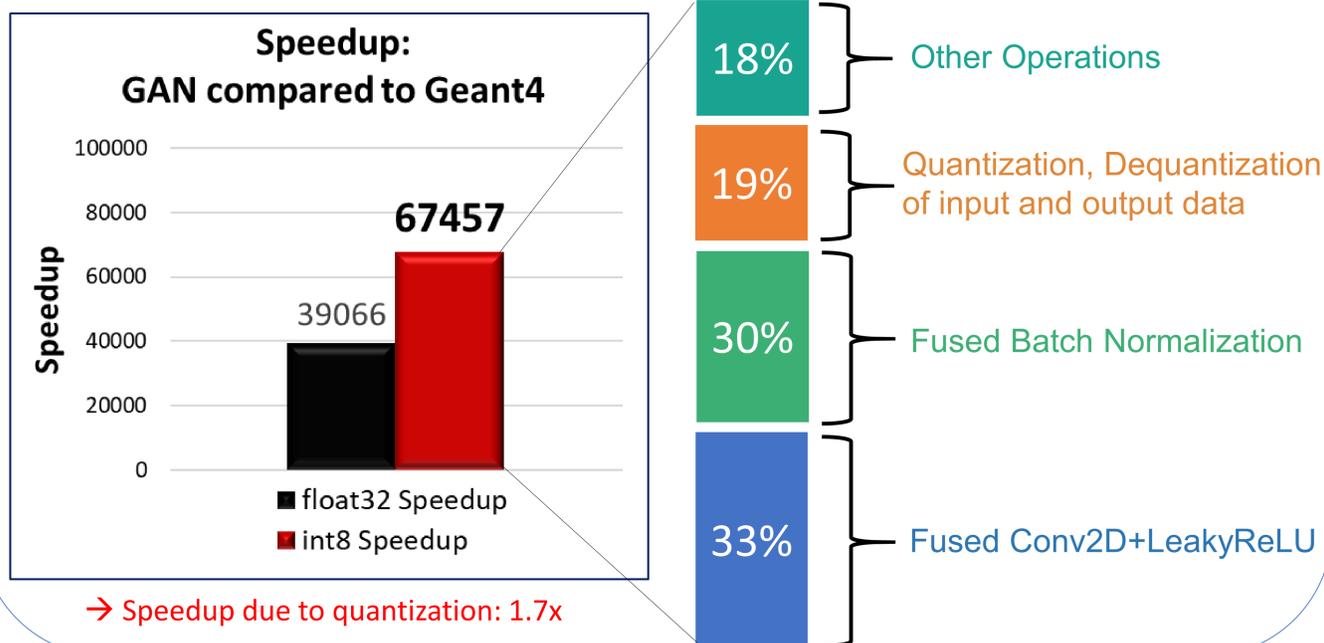
| Float32                     | Int8            |
|-----------------------------|-----------------|
| 4 byte                      | 1 byte          |
| Max Number: $3.4 * 10^{38}$ | Max Number: 255 |

- Geant4 → float32 GAN: 39000x speedup
- Geant4 → int8 GAN: **67000x** speedup
- Maintaining physics accuracy

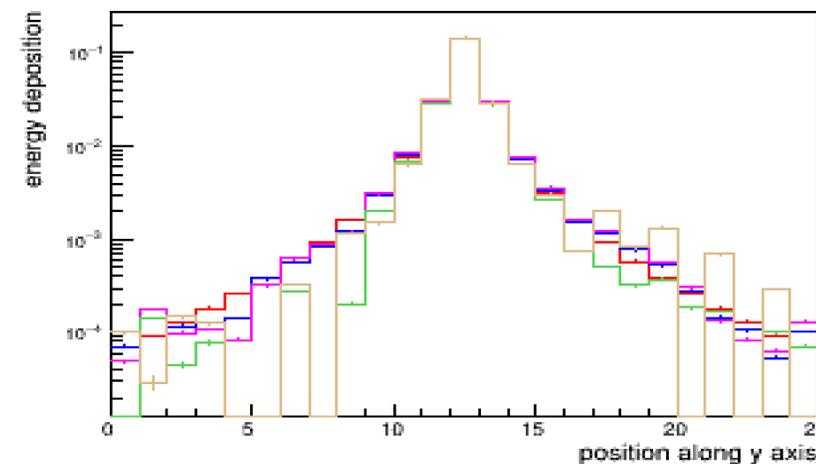
## Intel Low Precision Optimization Tool (iLoT)



## Computational Evaluation



## Physics Evaluation



| Model          | Uncertainty (Lower is better) |
|----------------|-------------------------------|
| float32        | 0.061                         |
| iLoT int8      | <b>0.053</b>                  |
| TFLite float16 | 0.253                         |
| TFLite int8    | 0.340                         |

→ Lower uncertainty  $\hat{=}$  better accuracy

