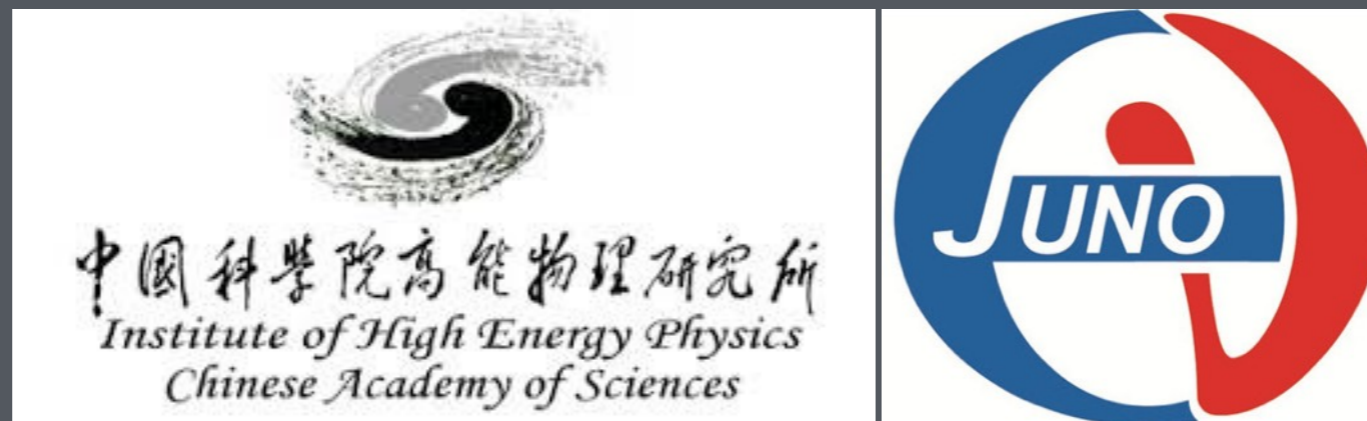# GPU APPLICATION IN JUNO

Wuming Luo
May 30th 2019

高 能 物 理 计 算 和 软 件 会 议
南 京 大 学
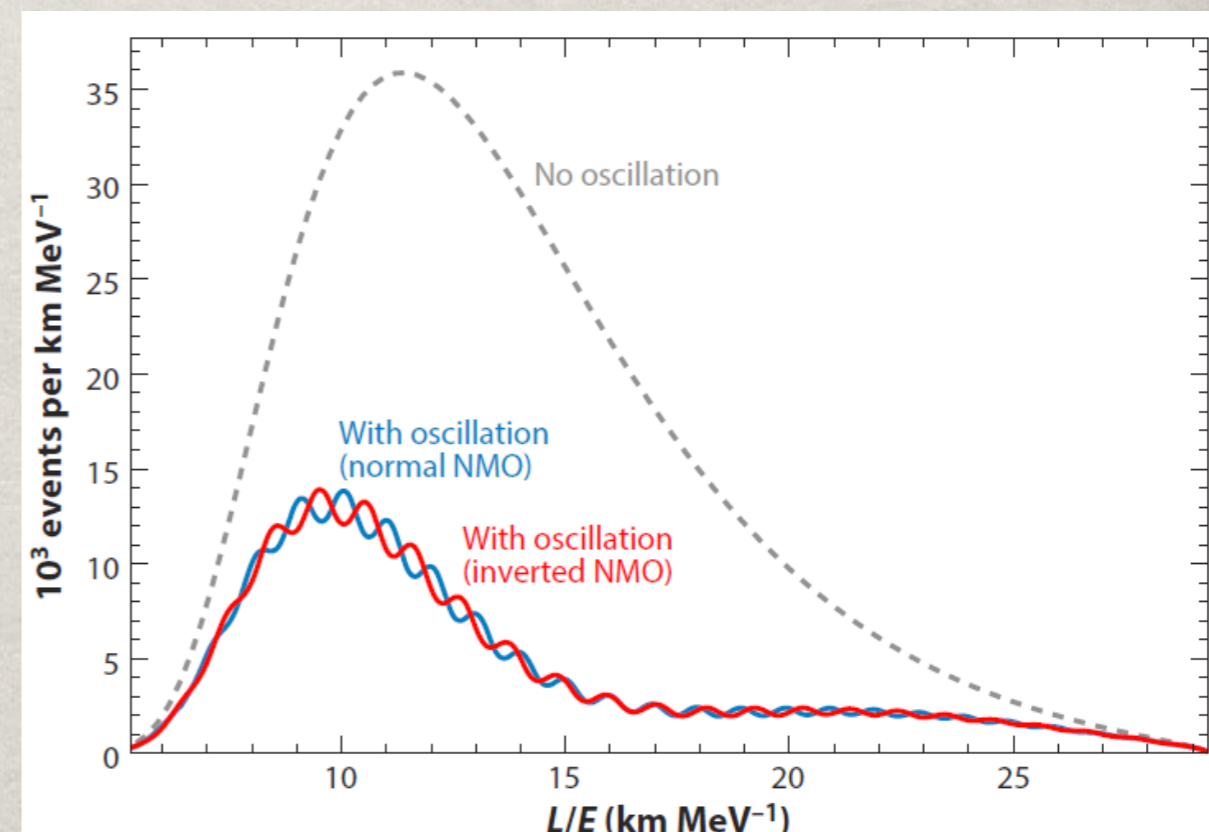
# OUTLINE

* Introduction to JUNO
* GPU vs CPU
* Applications
  * Vertex Reconstruction
  * Muon Simulation
  * Deep Learning
* Discussion
* Summary
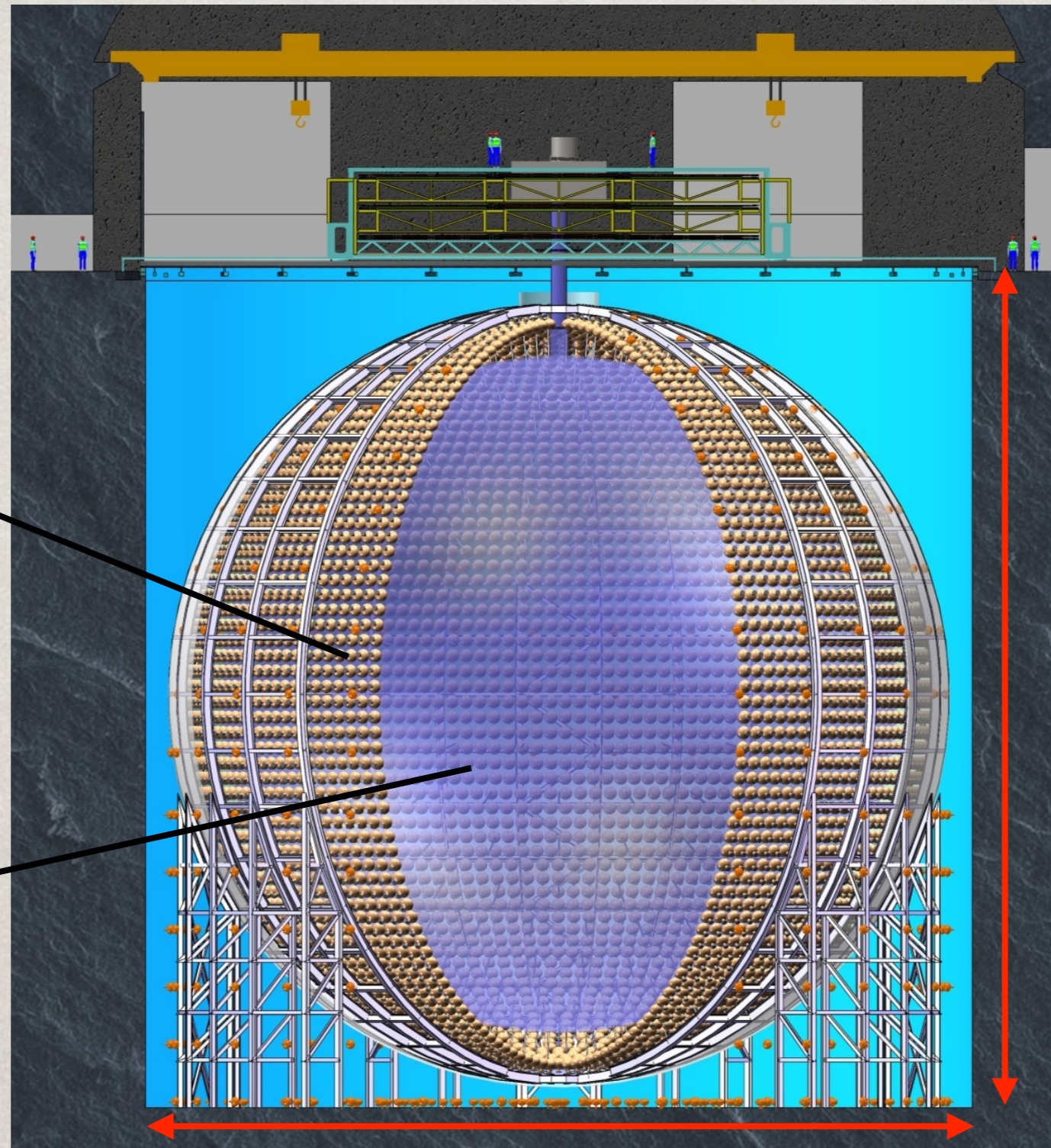
# JUNO

❋ Jiangmen Underground Neutrino Observatory(JUNO):

  ❋ Largest liquid scintillator detector (20 kton)

❋ Primary physics goals:

  ❋ Determine the neutrino mass hierarchy

  ❋ Measure neutrino oscillation parameters precisely

# DETECTOR



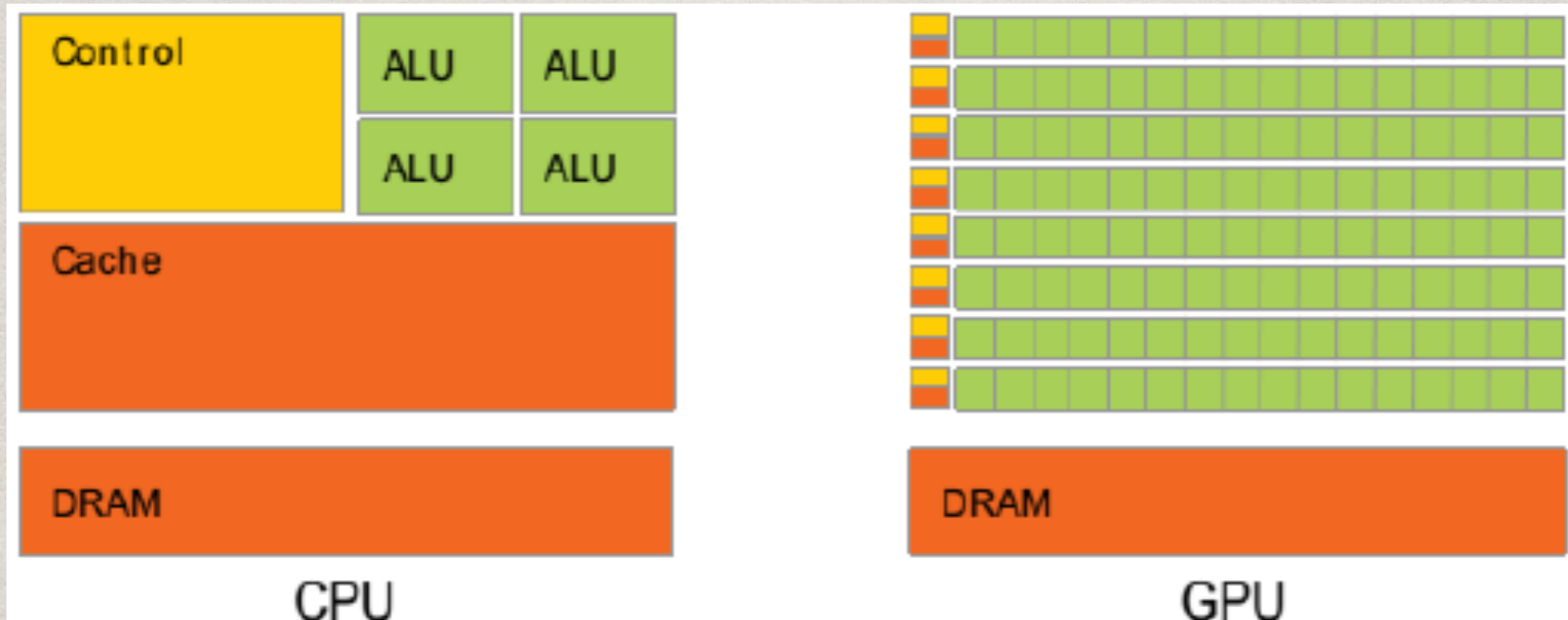Central Detector PMT
**~18,000 20" PMTs**
+ ~25,000 3" PMTs

**Liquid Scintillator**
20kton

Depth: 44m

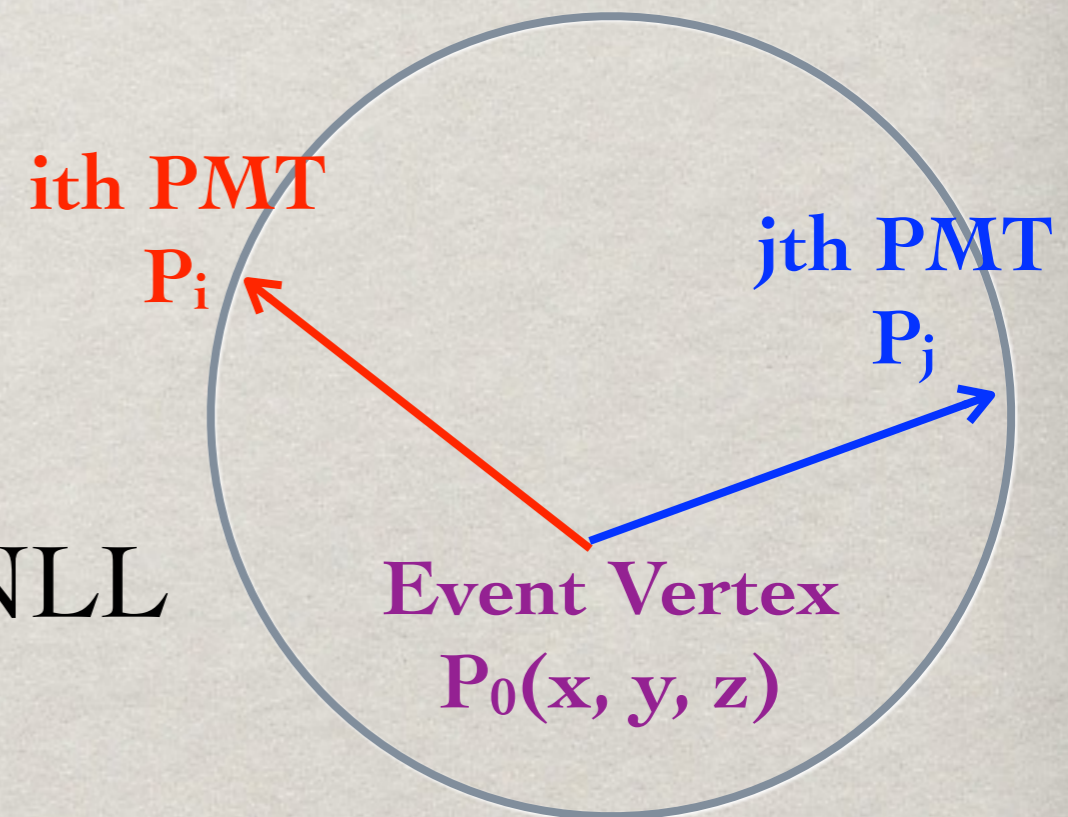$\phi$: 43.5m

Wuming Luo

4

# GPU VS CPU



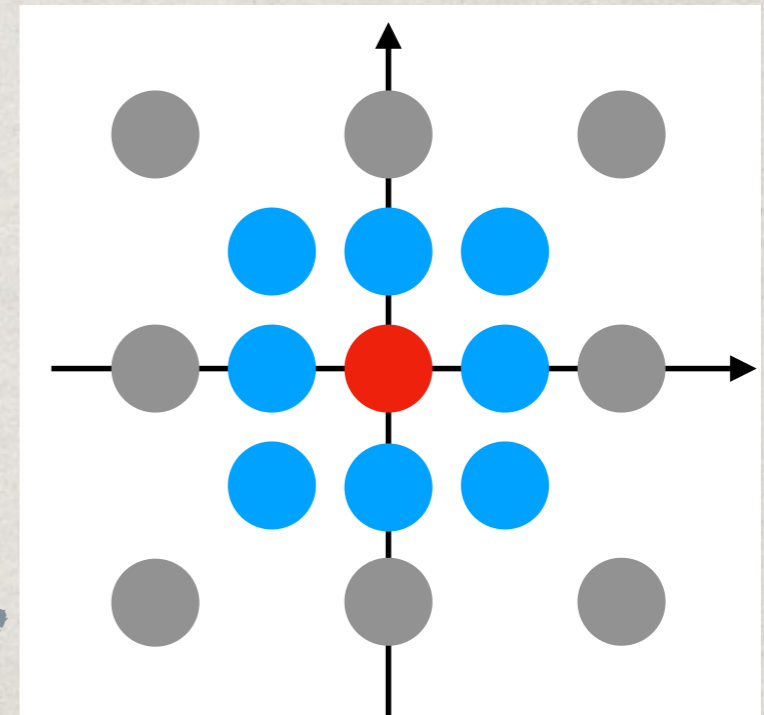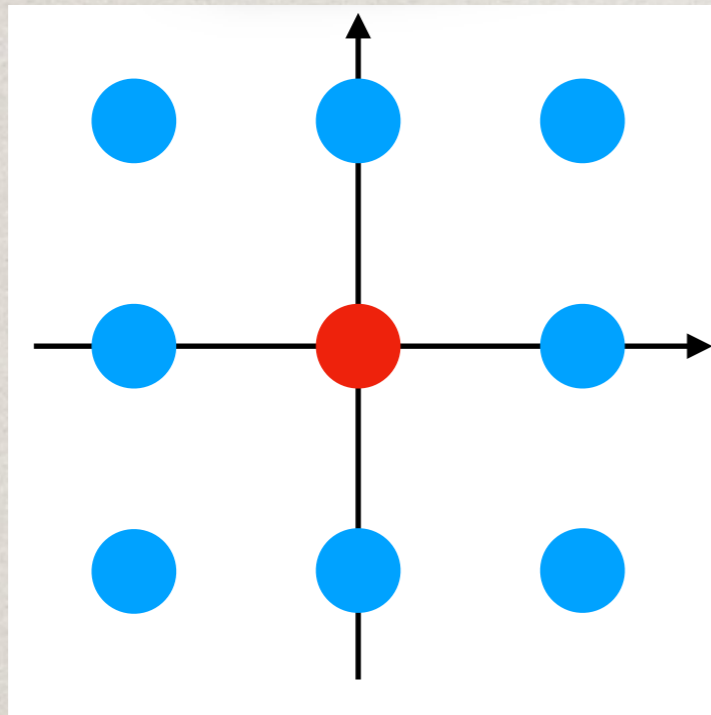Large Cache

Optimized for serial operations

Many cores

Built for parallel operations

# CASE 1: VERTEX RECONSTRUCTION

☀ Parameters to reconstruct: *x, y, z, t₀*
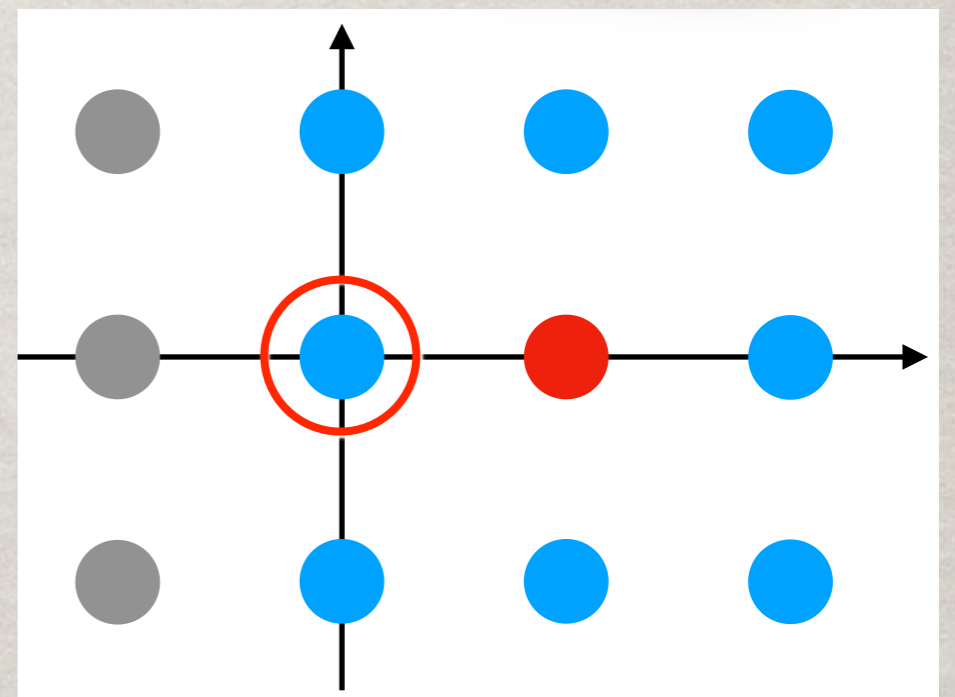
☀ Algorithm: $-\ln\mathscr{L} = -\sum \ln f_{res}(t_{i,res}) = -\sum \ln f_{res}(t_i - t_{i,tof} - t_0)$

  ☀ *$t_i$ : first hit time of ith fired PMT*

  ☀ *$t_{tof}$ : time of flight*

  ☀ *$t_0$ : event start time*

  ☀ *$f_{res}$ : pdf of residual time*

☀ Scan 4D grid to minimize the NLL

**ith PMT**
**$P_i$**

**jth PMT**
**$P_j$**

**Event Vertex**
**$P_0(x, y, z)$**

# GRID SEARCH — 2D

• Mini

• Nu

•

•

• Fo

• 24

•

*Schematic diagram*
*for 2-dim search*

```
if(Center is minimum){
    step /= 1/2
}
else{
    move to NEW center
}
```

mization alg                    earch

umber of gri

• 3 x-dim, 3

• 243 in total

or each grid

*diagram*

*search*  range from $3\times10^5$ (1 MeV) to $3\times10^6$ (10 MeV)

• Minimization algorithm: GridSearch, a 4-dim search

# PARALLELIZATION ON GPU

```
for(t){
    for(x){
        for(y){
            for(z){
                for(ith PMT){
                    calc. NLLi
                }
            ...
}
```
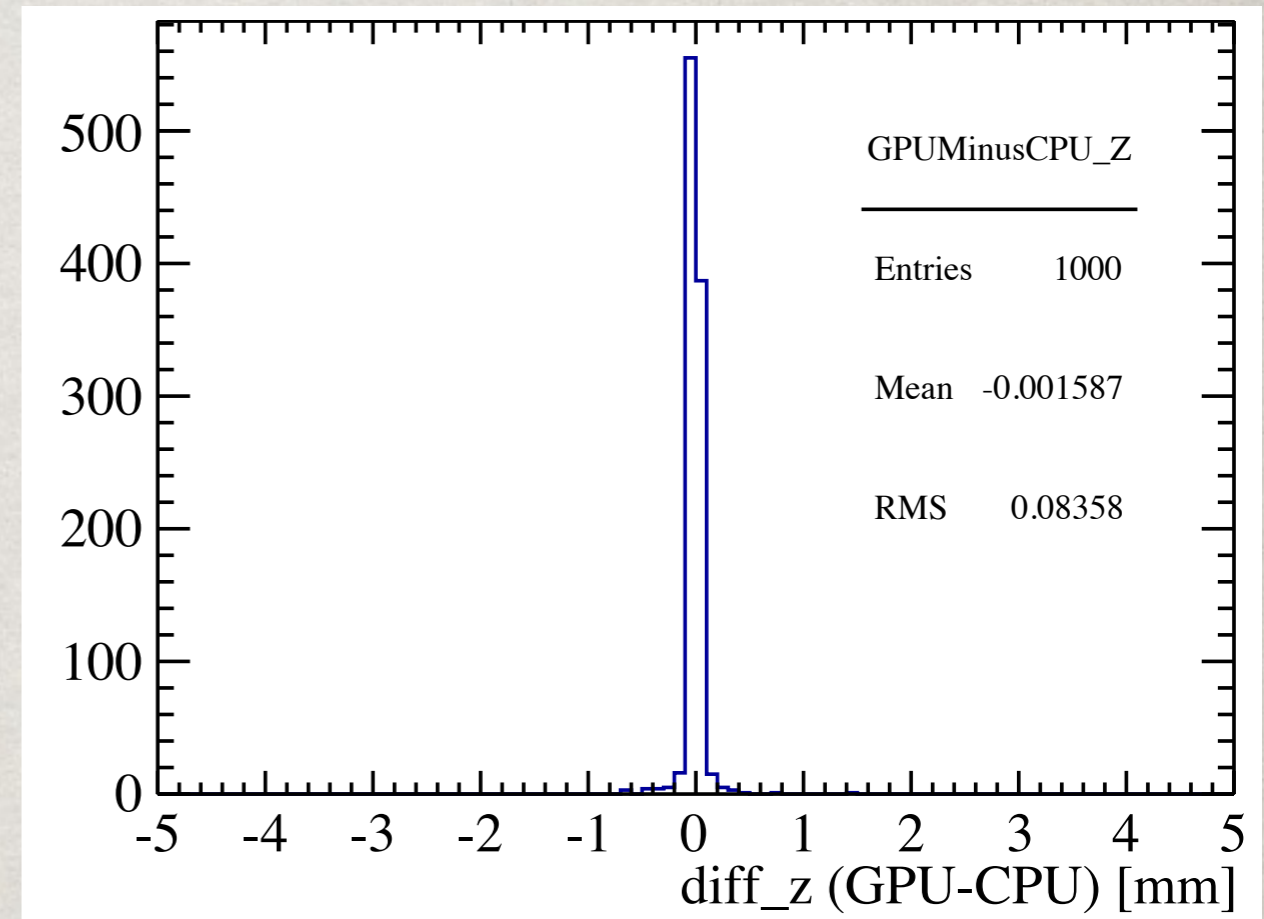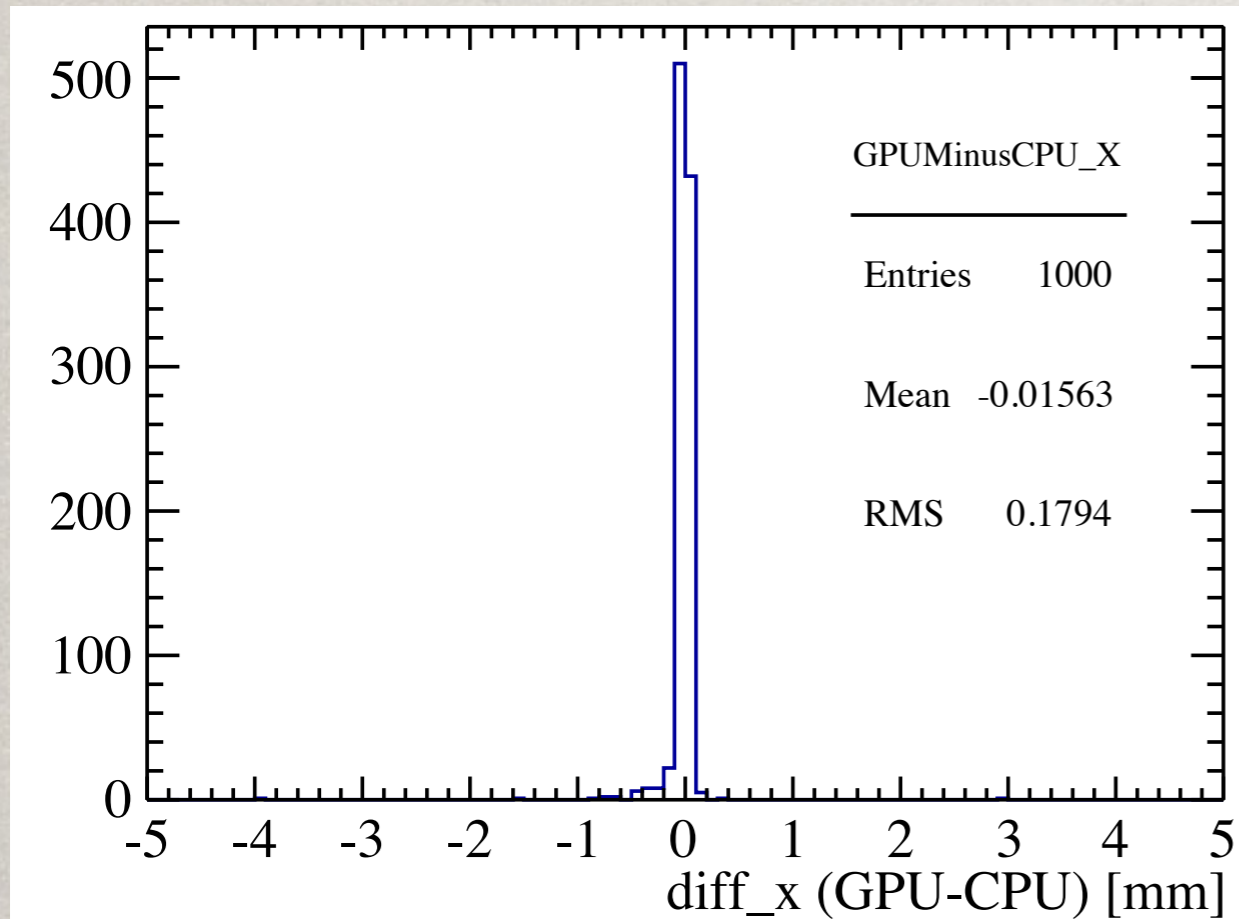ON CPU

- 4D Grid Search
- Number of loops: x-dim*y-dim*z-dim*t-dim*n_fired_PMTs = 3*3*3*9*1200/MeV = $3*10^5$/MeV
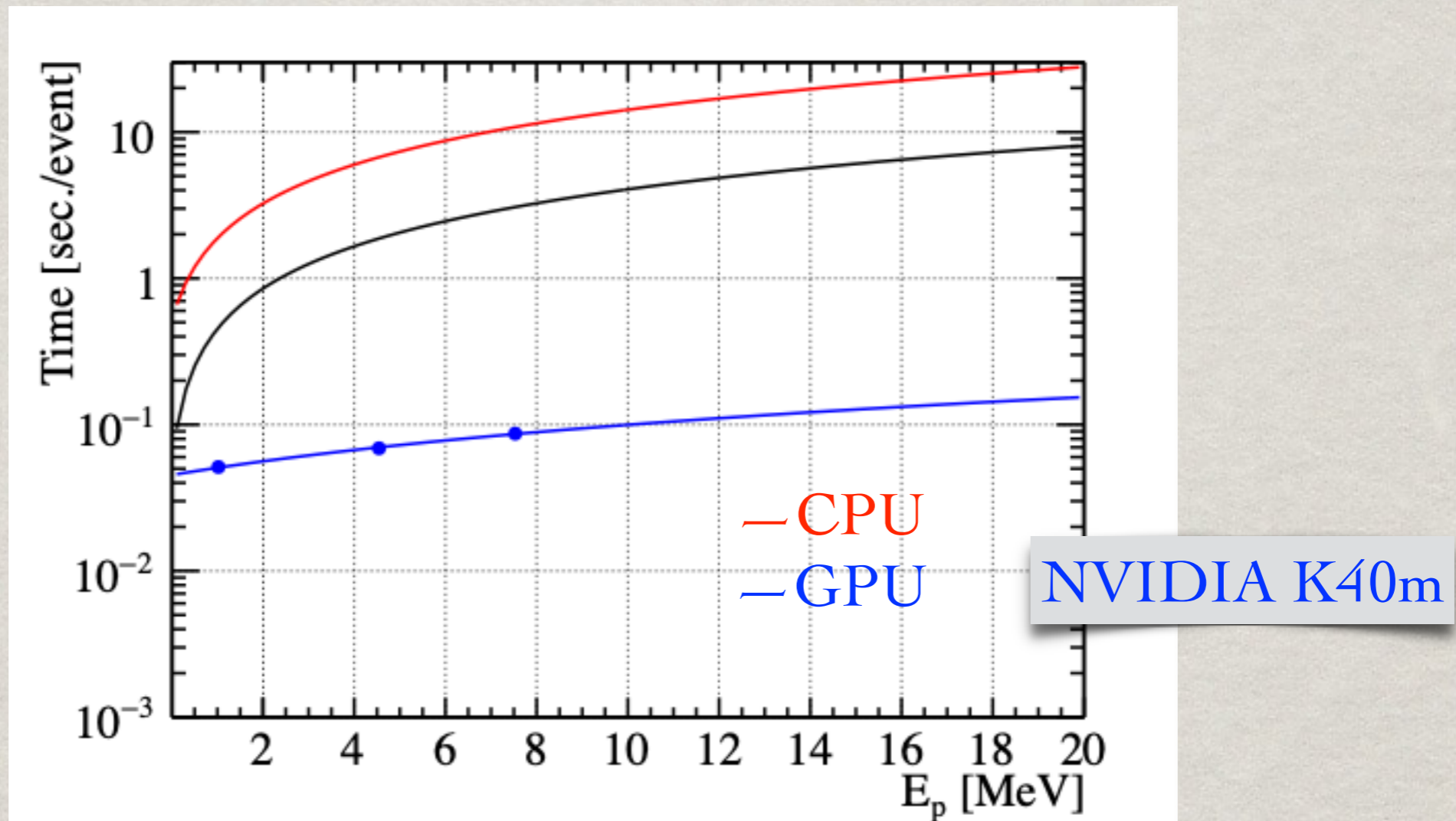- Parallelize the calculations on GPU

# VALIDATION



- GPU Rec was able to reproduce the CPU Rec results
- Tiny difference (<0.5mm), negligible w.r.t. vertex resolution (60mm)

# PERFORMANCE



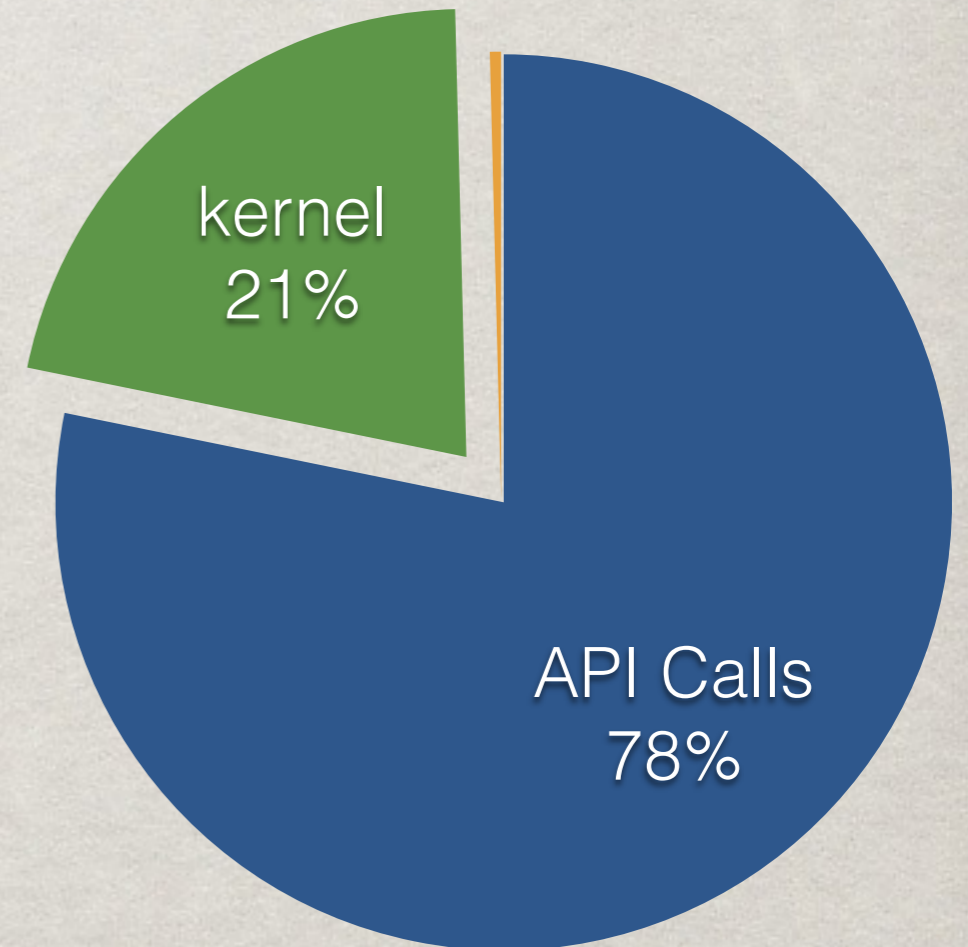|  | CPU | GPU | Ration: CPU/GPU |
|---|---|---|---|
| Time@1MeV(s) | 1.88 | 0.05 | ~40 |
| Time@10MeV(s) | 14.19 | 0.095 | ~150 |
| Gradiant | 1.37 | 0.005 | — |

# DISCUSSION

* Memory allocation and free, Synchronization etc… take up most of the time, room for future optimization

* Potential improvement with multiple GPUs

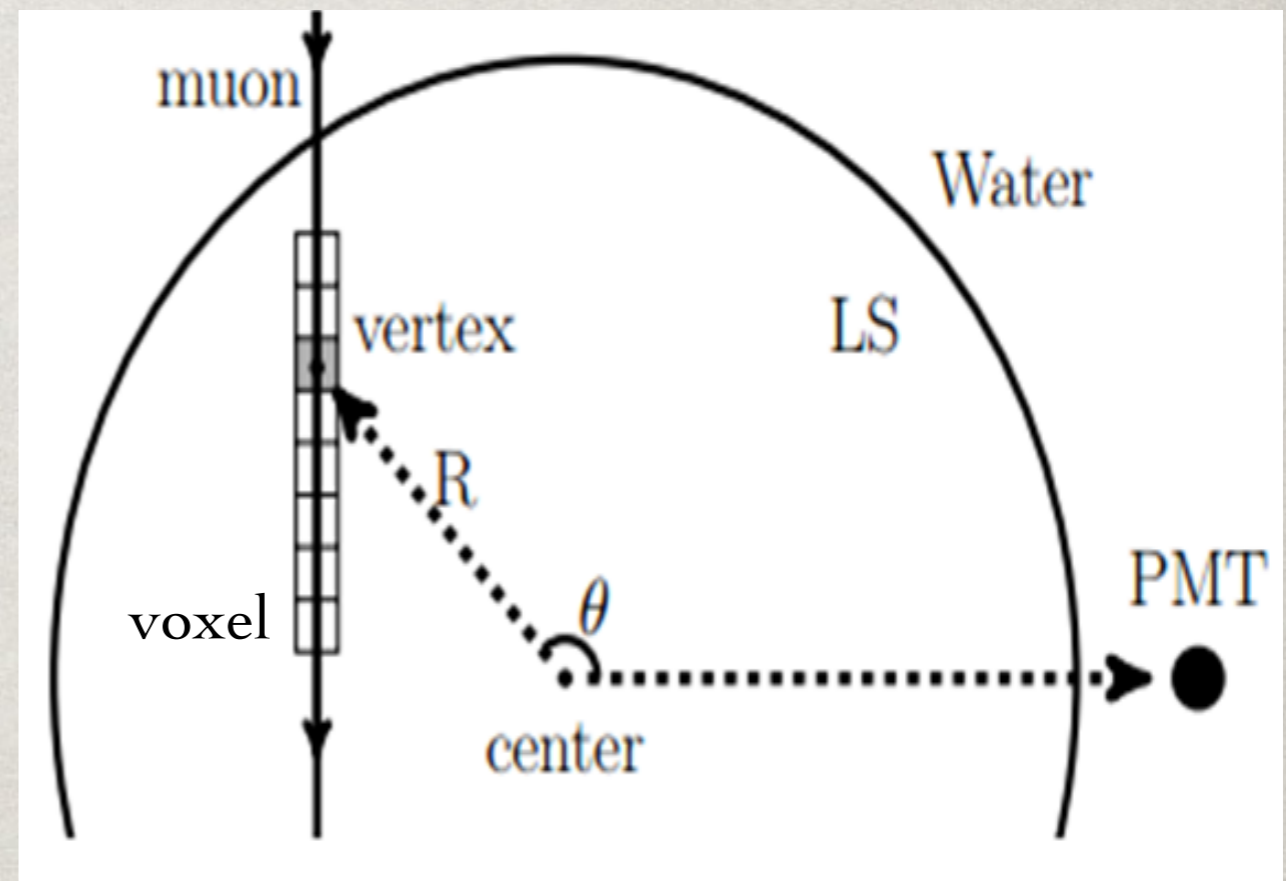* Instead of Grid Search, divide the detector ROI to tiny units and parallelize with GPU(s)

NVIDIA K40m

● API Calls    ● kernel
● data transfer

kernel
21%

API Calls
78%

* Simulate the number of photons (nPE) and the corresponding hit time($\{t_i\}$) collected by each PMT for a traversing Muon

* Voxel: segments along the muon track

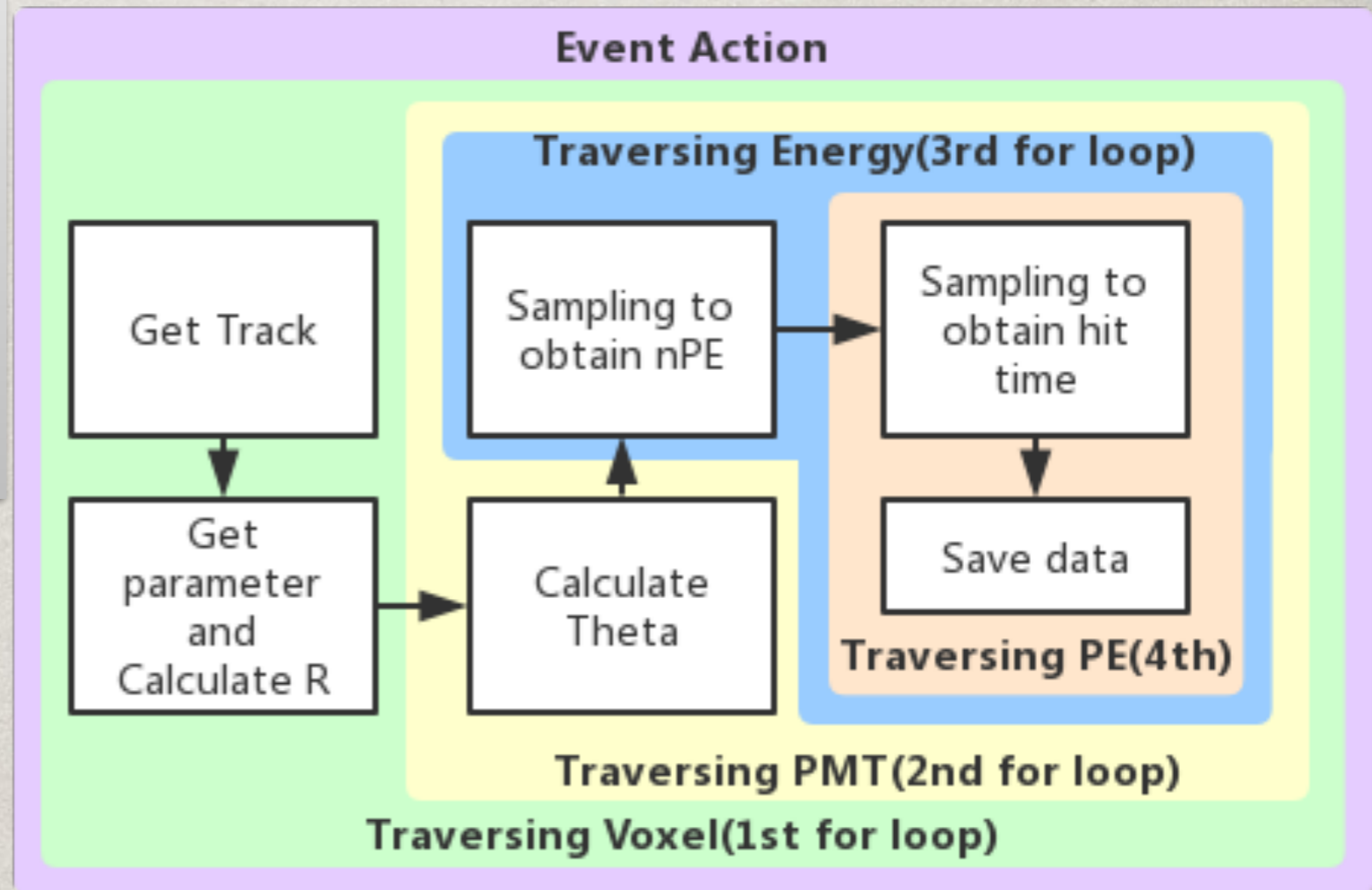* For fixed (R, $\theta$), sampling nPE and $\{t_i\}$ from templates

# COMPUTATION FLOW

for(R){      // Voxel loop
    for($\theta$){   // PMT loop
        for(E){   // E loop
            for(nPE){
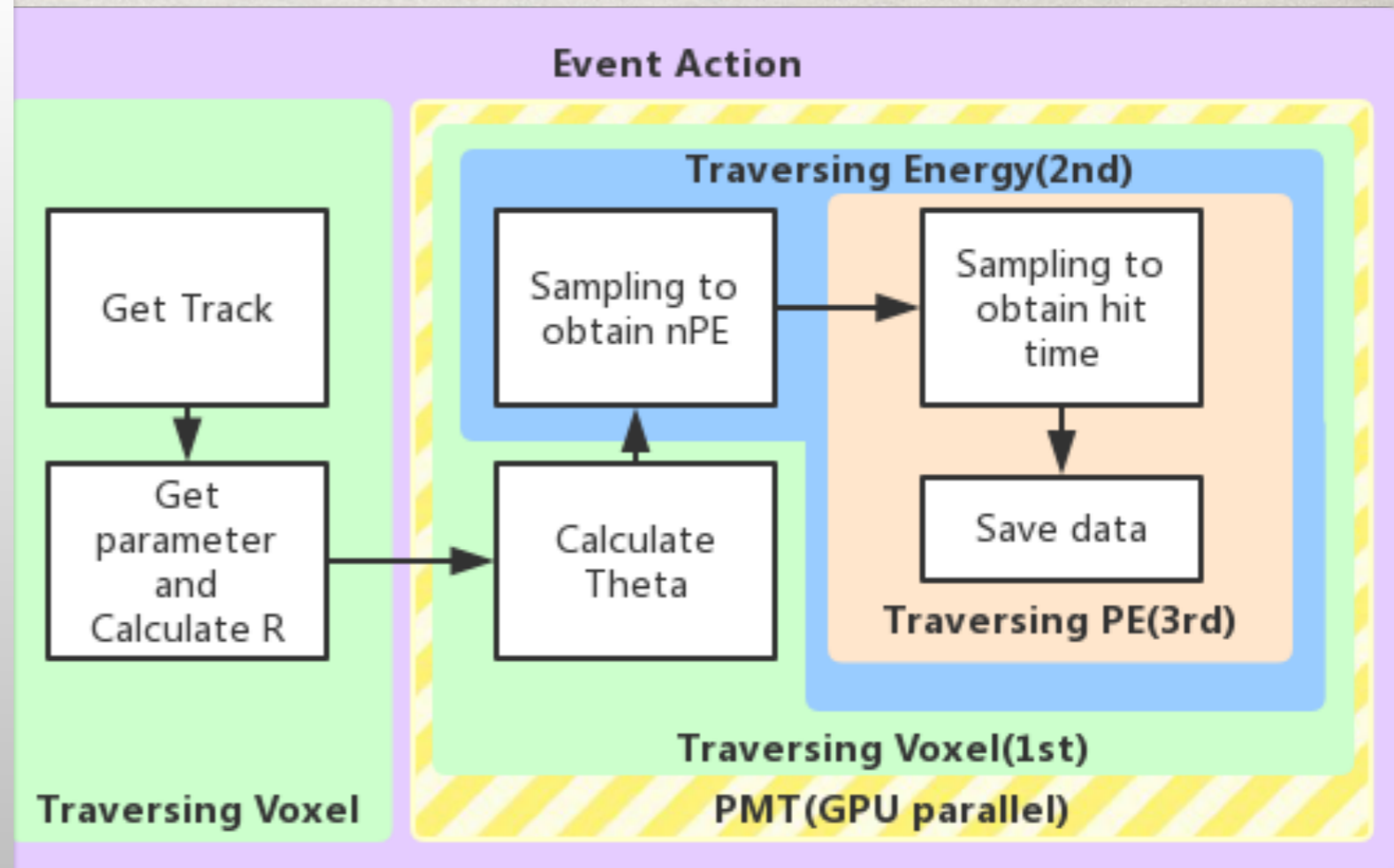                sample $t_i$
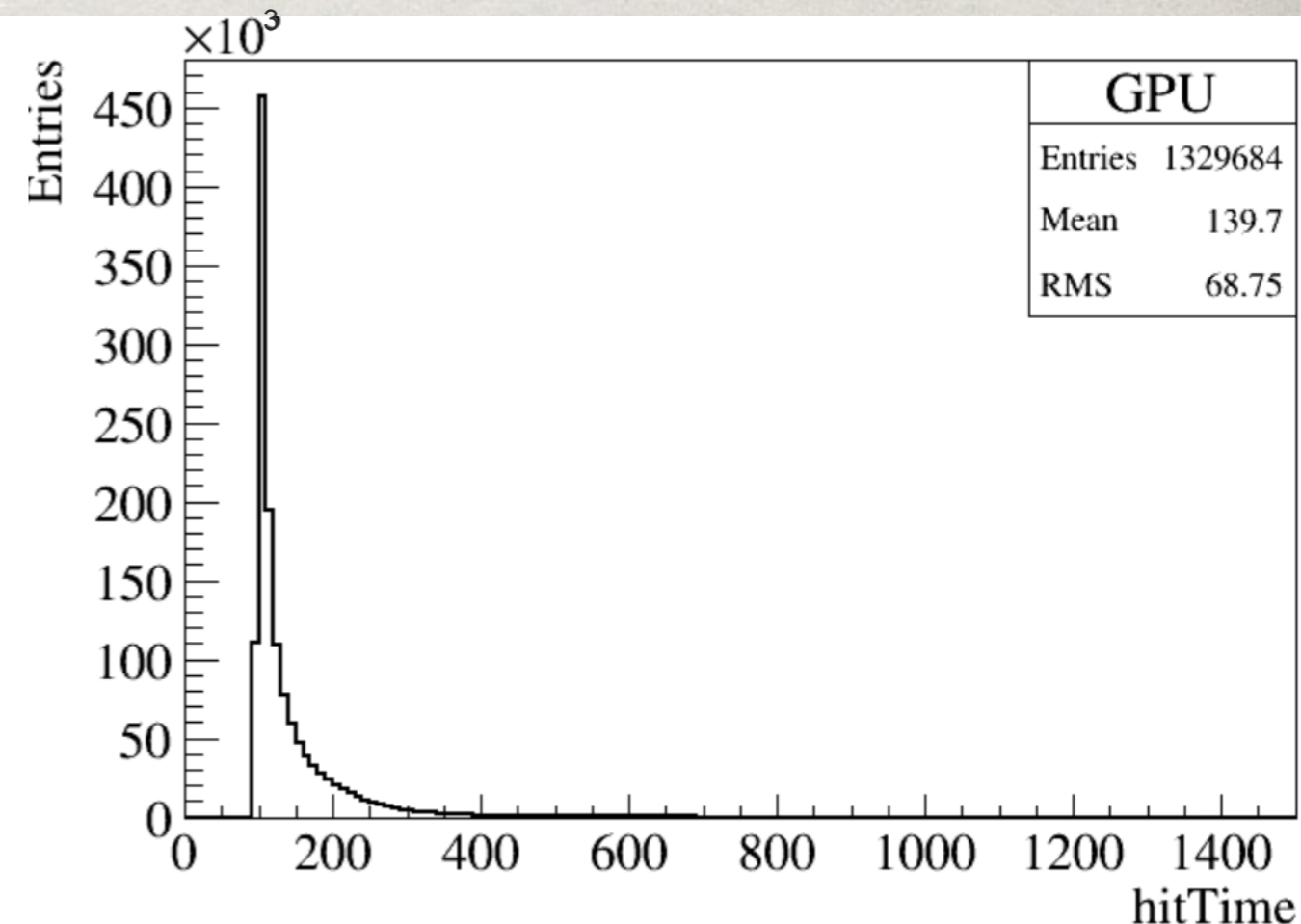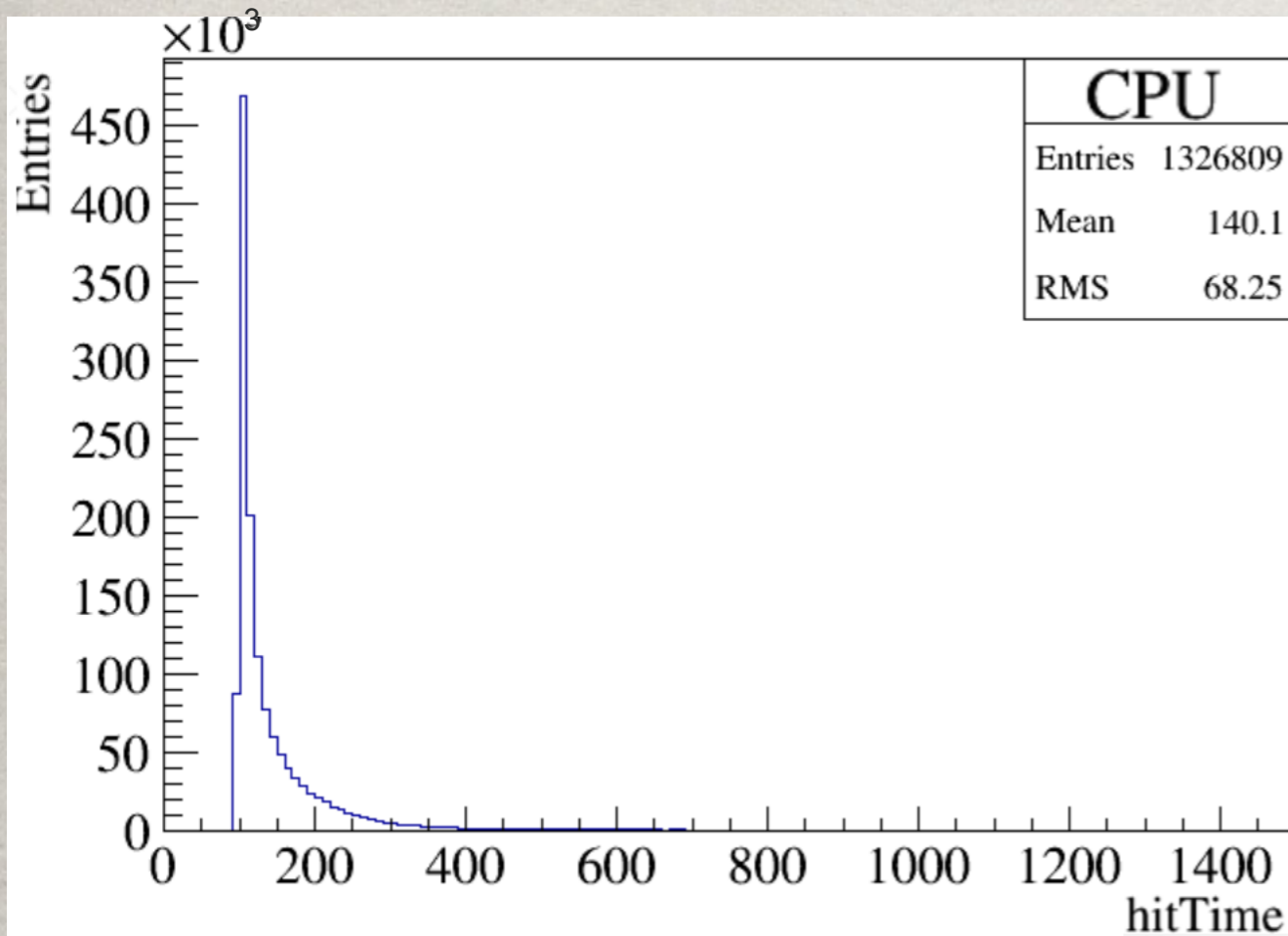
            }

    …
}

**ON CPU**

~18,000 PMTs

# COMPUTATION FLOW

```
for(θ){ // PMT loop
    for(R){ // Voxel loop
        for(E){ // E loop
            for(nPE){
                sample t_i
            }
        }
    ...
}
```

**Event Action**

**Traversing Voxel**

Get Track → Get parameter and Calculate R → Calculate Theta

**PMT(GPU parallel)**

**Traversing Voxel(1st)**

**Traversing Energy(2nd)**

Sampling to obtain nPE → Sampling to obtain hit time
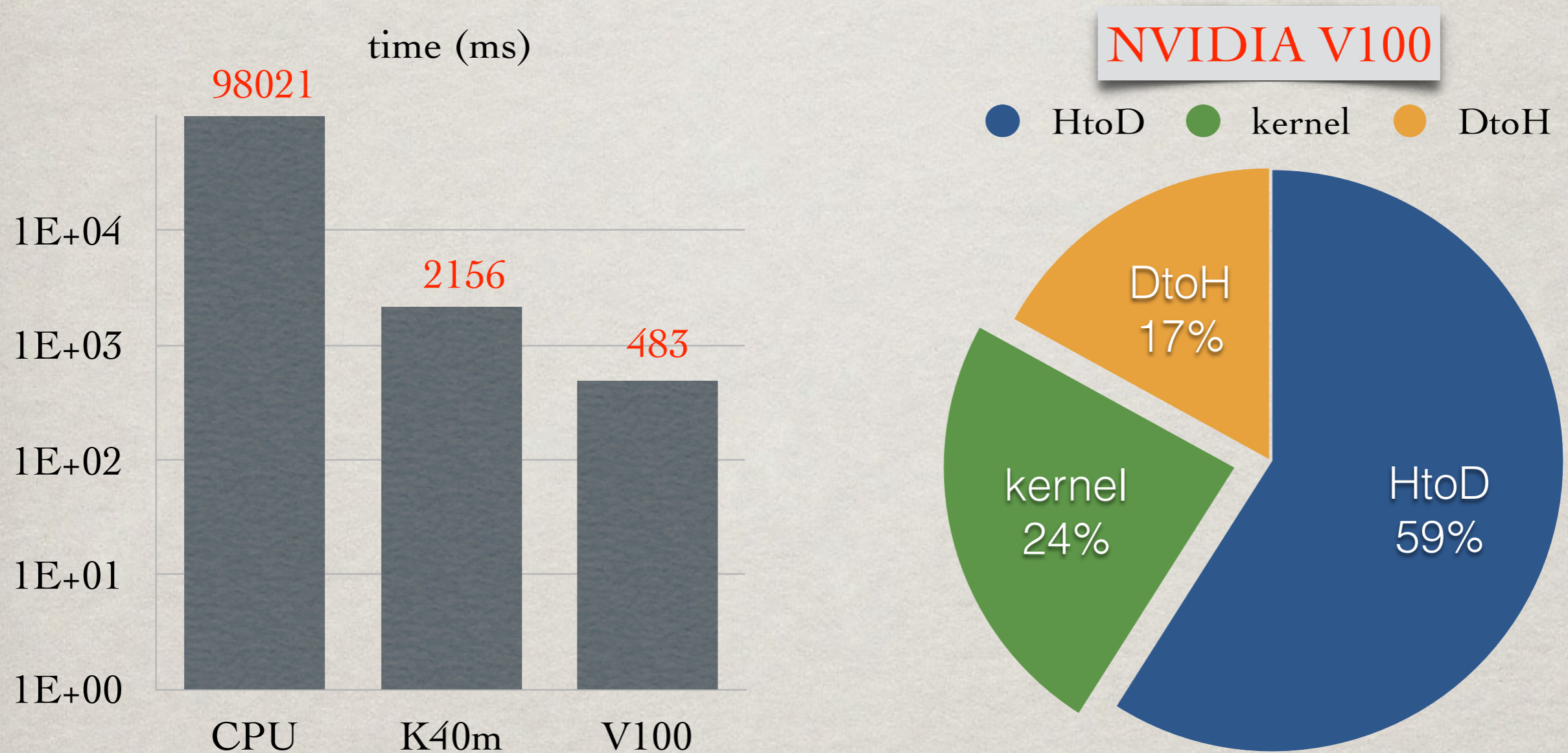
**Traversing PE(3rd)**

Save data

✳ Switch the Voxel loop and PMT loop levels

✳ Parallelize the PMT loop with GPU

# VALIDATION



* GPU Sim was able to reproduce the CPU Sim results
* Negligible difference

# PERFORMANCE

time (ms)

98021

2156

483

| | | |
|---|---|---|
| 1E+04 | | |
| 1E+03 | | |
| 1E+02 | | |
| 1E+01 | | |
| 1E+00 | | |
| CPU | K40m | V100 |

NVIDIA V100

● HtoD   ● kernel   ● DtoH

DtoH
17%

kernel
24%

HtoD
59%

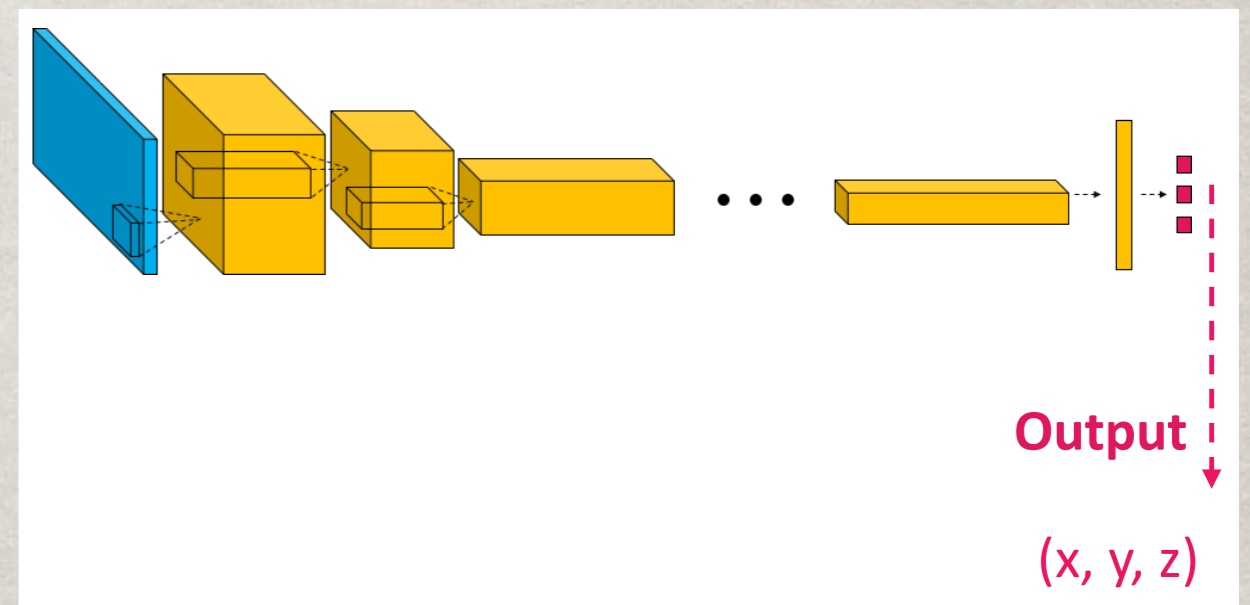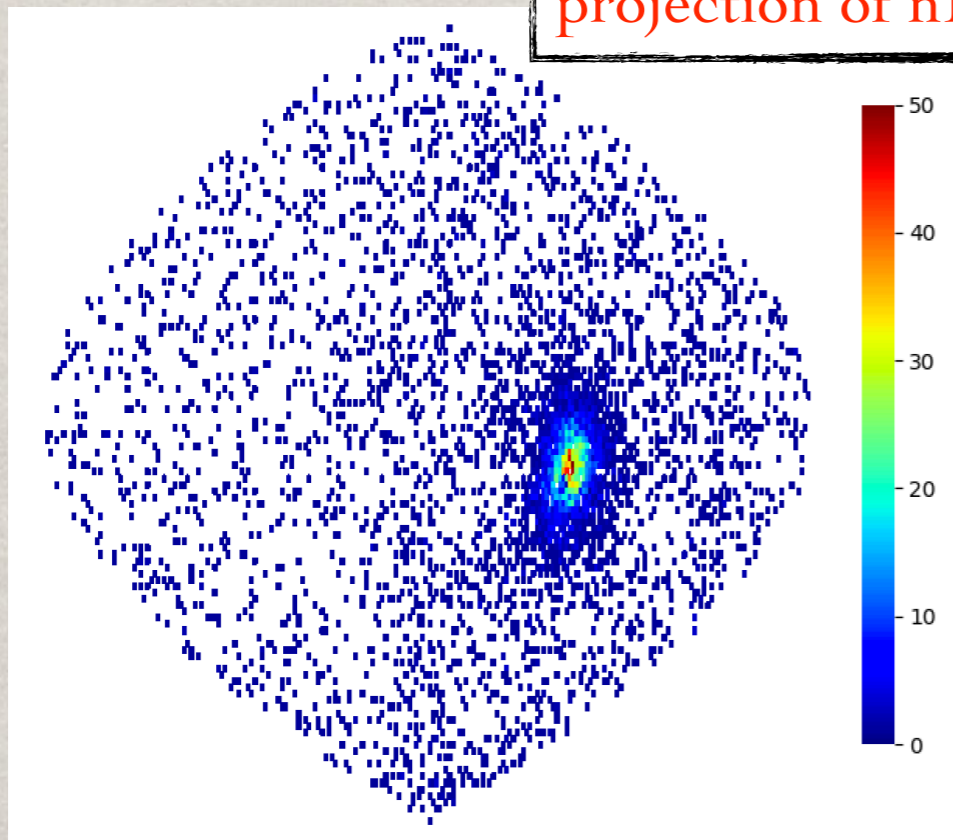※ O($10^2$) improvement with V100

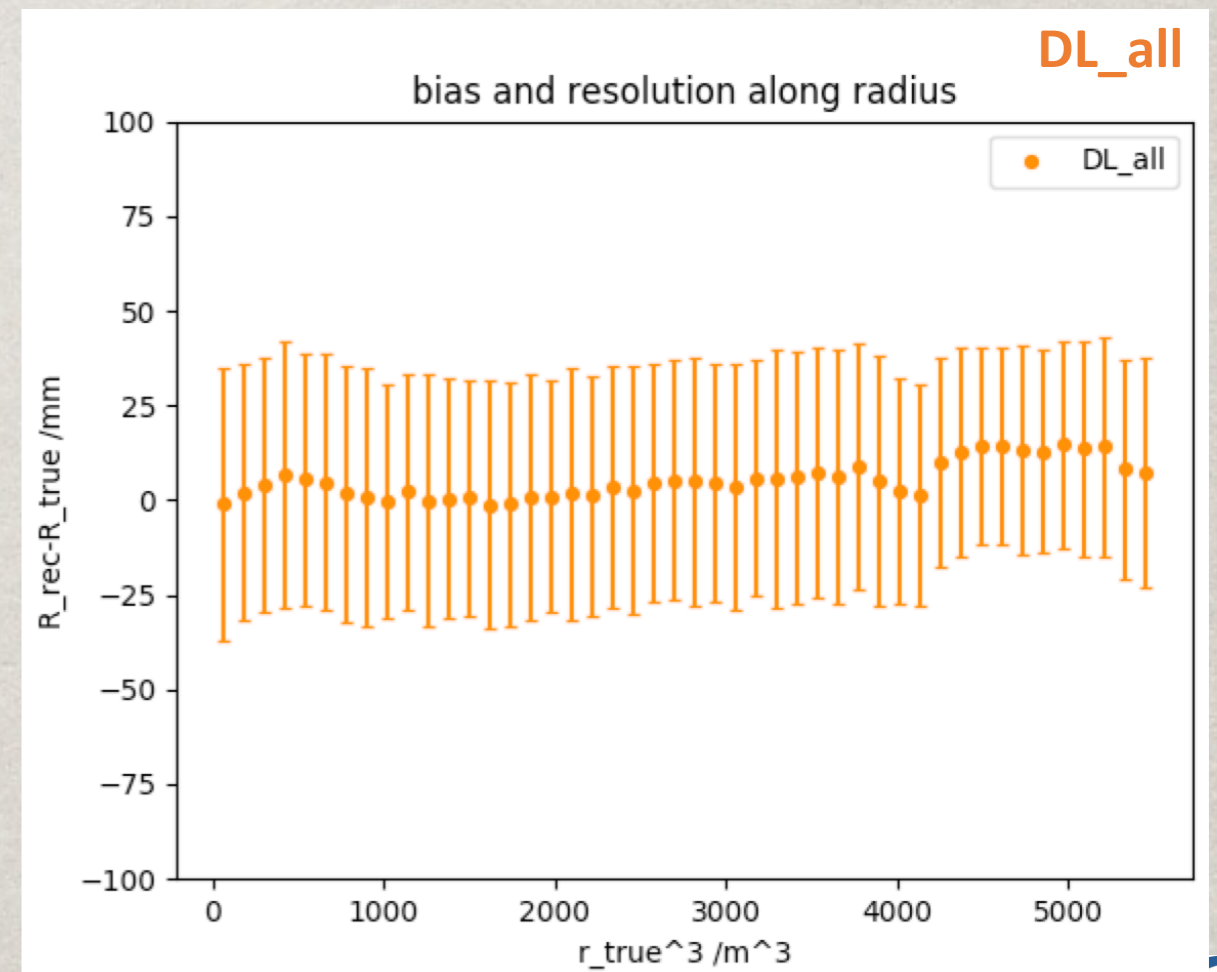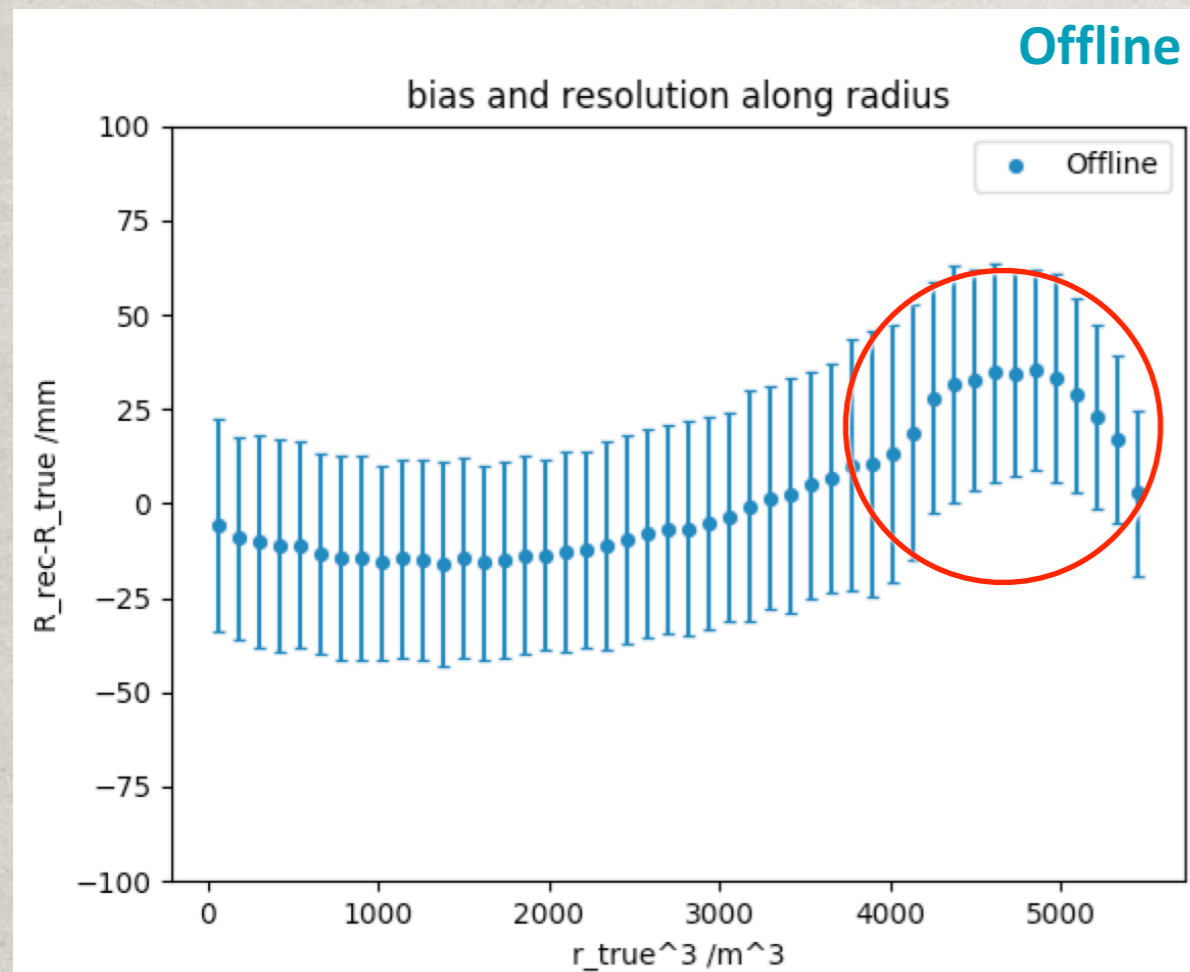※ Future optimization: data transfer, more levels, multi-GPUs,

JUNO

# CASE 3: DEEP LEARNING

⚜ GPU is widely used for DL

⚜ Try Vertex Reconstruction with CNN in JUNO

⚜ Input: hit time $\{t_i\}$, numbe            ons $\{nPE_i\}$

⚜ Output: event vertex $(x, y$

projection of nPE

Output

(x, y, z)

# PERFORMANCE
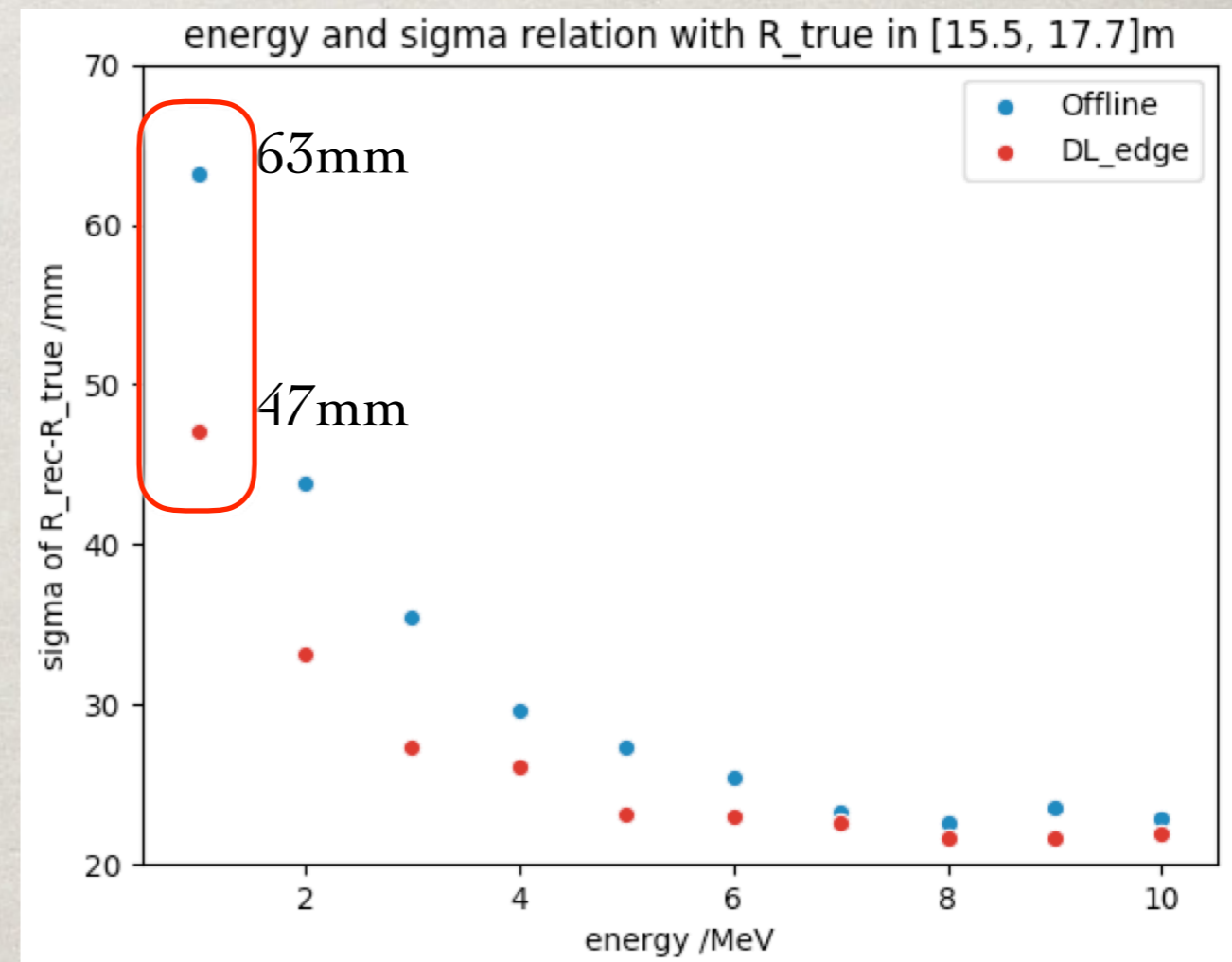
- Smaller Bias for DL w.r.t. traditional method
- Better uniformity for DL w.r.t. traditional method

# PERFORMANCE

Resolution vs E

- For events near the detector edge, total refraction complicates the optical model

- Better resolution for DL w.r.t. traditional method using a specialized model



energy and sigma relation with R_true in [15.5, 17.7]m

63mm

47mm

Offline
DL_edge

sigma of R_rec-R_true /mm

energy /MeV

# DISCUSSION FOR DL

* Pros:
  * fast speed, energy independent
  * avoid the complex optical model
* Cons:
  * rely heavily on GOOD Monte Carlo simulation
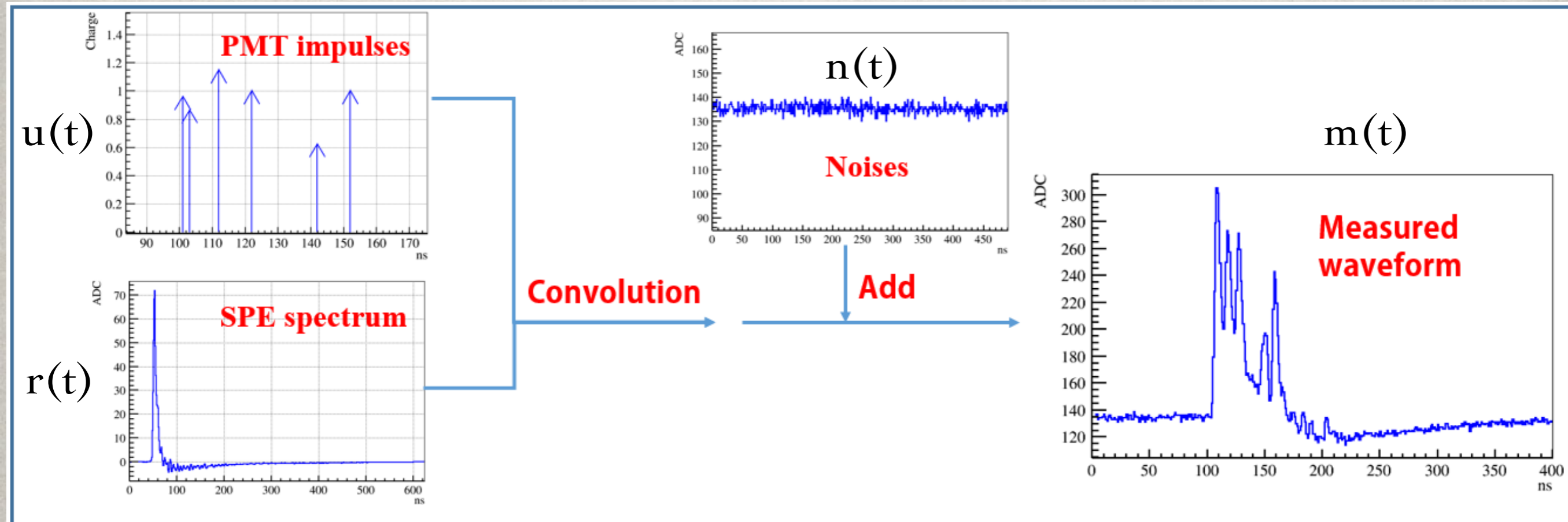* Training samples
  * MC: large statistics, might be different w.r.t. real data
  * Calibration data: close to real data, limited stats.
* Possible solutions?

# PMT WAVEFORM REC



- $m(t) = s(t) + n(t) = r(t)^{*}u(t) + n(t)$
- We need to reconstruct $\{t_j\}$ and $\{charge_j\}$ or ideally $\{nPE_j\}$

# DL FOR WAVEFORM REC?

* FADC raw waveform —> Time series
* We know roughly what the feature looks like —> sPE response template
* We want to know $\{t_j, Q_j(nPE_j)\}$ for all pulses
* We have PMT testing data —> real waveform
  * Issue: unsupervised, real labels unknown
* Analogies? Voice recognition? Suggestions?
* Try to answer simpler questions:
  * Q1: what is the first hit time?
  * Q 2: classify waveform to [0, 1, ≥2]PE three categories

# SUMMARY

- JUNO has ~$O(10^5)$ PMTs, perfectly suitable for utilizing GPU
- Showed a few simple applications of GPU in JUNO
  - Vertex reconstruction, Muon simulation, Deep Learning
  - Large room for further improvements
- Could be used in other aspects of JUNO
- Huge potential for experiments with lots of PMTs

# BACKUP

# TOOLS

* CUDA
* Thrust
* TensorFlow

| | multi-processors | CUDA cores | ram(GB) |
|---|---|---|---|
| K40m | 15 | 2880 | 12 |
| V100 | 80 | 5120 | 32 |