# Energy reconstruction with machine learning techniques in JUNO: aggregated features approach

**Arsenii Gavrikov**[1,2], Yury Malyshkin[2], Fedor Ratnikov[1]

[1]HSE University, Moscow, Russia

[2]Joint Institute for Nuclear Research, Dubna, Russia

Moscow International School of Physics 2022, 24 July – 2 August 2022

# Introduction to the JUNO experiment

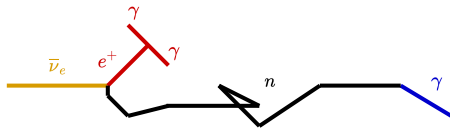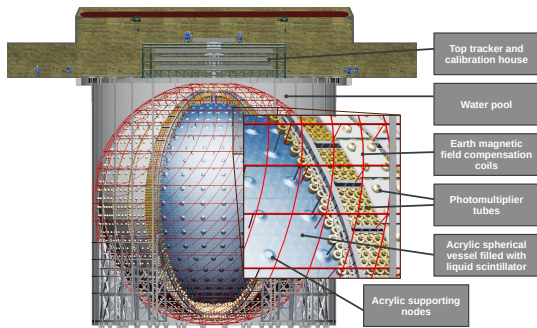1. Jiangmen Underground Neutrino Observatory:
   - multipurpose experiment
   - 53 km away from 8 reactor cores in China
   - data taking expected in $\sim$2023
   - JUNO Collaboration:
     - 77 institutions
     - 697 collaborators

2. The main goals of JUNO:
   - neutrino mass ordering ($3\sigma$ in 6 years)
   - precise measure of oscillation parameters $\sin^2 \theta_{12}, \Delta m_{21}^2, \Delta m_{31}^2$;

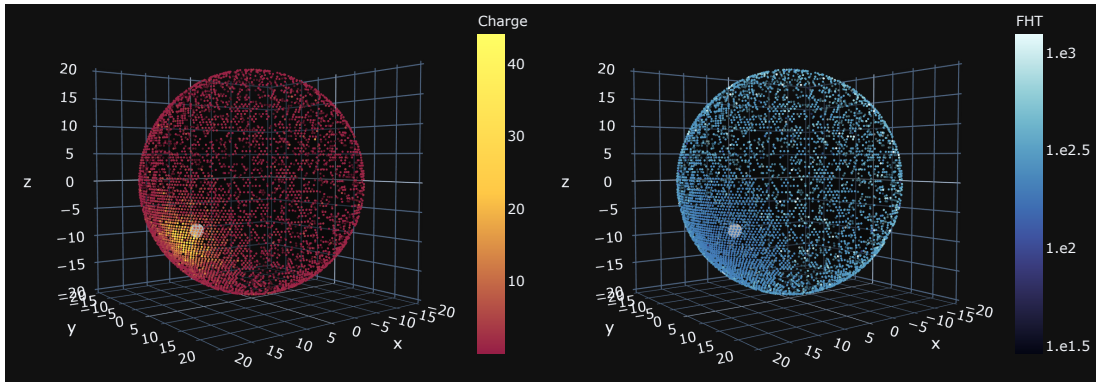3. The Central Detector:
   - detection channel: $\overline{\nu}_e + p \rightarrow e^+ + n$;
   - deposited energy converts to optical light
   - the largest liquid scintillator detector: 20 kt
   - 77.9% photo-coverage: 18k 20", 26k 3" photo-multiplier tubes (PMTs)



Top tracker and calibration house

Water pool

Earth magnetic field compensation coils

Photomultiplier tubes

Acrylic spherical vessel filled with liquid scintillator

Acrylic supporting nodes

# Machine Learning (ML) in HEP

- ML methods are used at all levels of data processing in many HEP experiments:
  - signal/background discrimination
  - event selection in the trigger
  - event simulation
  - anomaly detection
  - identification, etc.

- Why is ML useful for HEP?
  - **Faster**. More precisely, with proper training
  - **Adequate** for many purposes simultaneously: event simulation, analysis, reconstruction, identification, etc.
  - **GPU friendly** by construction, which is important for big data processing

- Machine-learning algorithms use statistics to find patterns in massive amounts of data

- Our task is a supervised learning problem (regression)

# Problem statement



An example of a positron event with deposited energy ∼**6 MeV**. The grey sphere — the primary vertex.

🟥 Charge at PMT      🟦 First Hit Time (FHT) at PMT

**We want to reconstruct**:

Deposited energy $E_{\text{dep}}$ with resolution 3% @ 1 MeV

## Datasets

- Two datasets: for training and for testing
- generated by the Monte Carlo method
- full detector and electronics simulation
- using the official JUNO software

**Data description:**

1. positron events
2. uniformly spread in the volume of the central detector
3. $E_{\text{kin}} \in [0, 10]$ MeV. $E_{\text{dep}} = E_{\text{kin}} + 1.022$ MeV

- **Training dataset:**
    4. **5 million** events
    5. uniformly distributed in kinetic energy $E_{\text{kin}}$

- **Testing dataset:**
    4. subsets with discrete kinetic energies:
    5. 0, 0.1, 0.3, 0.6, 1, 2, ..., 10 [MeV]
    6. $\sum = $ **1.4 million** events: each subset contains 100k

# Aggregated features

We use aggregated information from the whole array of PMTs as features for models:

1. `AccumCharge` — the accumulated charge on fired PMTs

2. `nPMTs` — the total number of fired PMTs

3. Coordinates of the center of charge:

$$(x_{cc},\ y_{cc},\ z_{cc}) = \vec{r}_{cc} = \frac{\sum_{i=1}^{N_{PMTs}} \vec{r}_{PMT_i} \cdot n_{p.e.,i}}{\sum_{i=1}^{N_{PMTs}} n_{p.e.,i}}$$

   and its radial component: $R_{cc} = |\vec{r}_{cc}|$

4. Coordinates of the center of FHT:

$$(x_{cht},\ y_{cht},\ z_{cht}) = \vec{r}_{cht} = \frac{1}{\sum_{i=1}^{N_{PMTs}} \frac{1}{t_{ht,i}+c}} \sum_{i=1}^{N_{PMTs}} \frac{\vec{r}_{PMT_i}}{t_{ht,i}+c},$$

   and its radial component: $R_{cht} = |\vec{r}_{cht}|$

5. $\gamma_z^{cc} = \frac{z_{cc}}{\sqrt{x_{cc}^2+y_{cc}^2}}$

6. $\gamma_y^{cc} = \frac{y_{cc}}{\sqrt{x_{cc}^2+z_{cc}^2}}$

7. $\gamma_x^{cc} = \frac{x_{cc}}{\sqrt{z_{cc}^2+y_{cc}^2}}$

8. $\theta_{cc} = \arctan \frac{\sqrt{x_{cc}^2+y_{cc}^2}}{z_{cc}}$

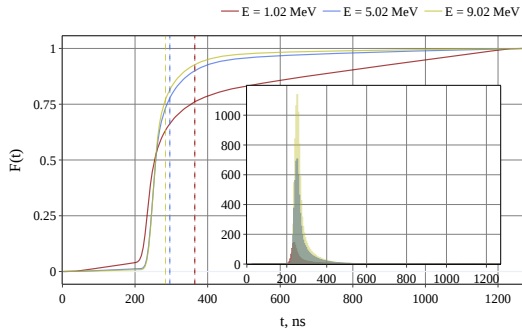9. $\phi_{cc} = \arctan \frac{y_{cc}}{x_{cc}}$

10. $J_{cc} = R_{cc}^2 \cdot \sin\theta_{cc}$

11. $\rho_{cc} = \sqrt{x_{cc}^2 + y_{cc}^2}$

12. with 7 similar features for the components of the center of FHT

# Aggregated features

**13** Percentiles of FHT and charge distributions:
- $\{ht_{2\%}, ht_{5\%}, ht_{10\%}, ht_{15\%}, ..., ht_{90\%}, ht_{95\%}\}$
- $\{pe_{2\%}, pe_{5\%}, pe_{10\%}, pe_{15\%}, ..., pe_{90\%}, pe_{95\%}\}$

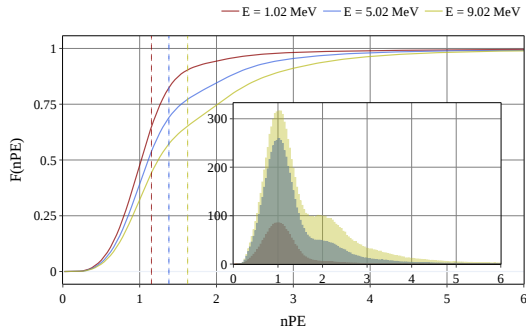**14** Differences between percentiles for FHT:
- $\{ht_{5\%-2\%}, ht_{10\%-5\%}, ..., ht_{95\%-90\%}\}$

**15** Moments for FHT and charge distributions:
- $\{ht_{mean}, ht_{std}, ht_{skew}, ht_{kurtosis}\}$
- $\{pe_{mean}, pe_{std}, pe_{skew}, pe_{kurtosis}\}$



CDFs and PDFs for FHT (left) and charge (right) distributions. $R \simeq 0$ m, $E_{kin}$ varied. Dashes lines show mean values.

# Models description: BDT

A Decision Tree (DT) takes a set of input features and splits input data recursively based on those features.

**Boosted Decision Trees** (BDT):

- Ensemble model
- DT as base algorithm
- DTs in BDT are trained sequentially
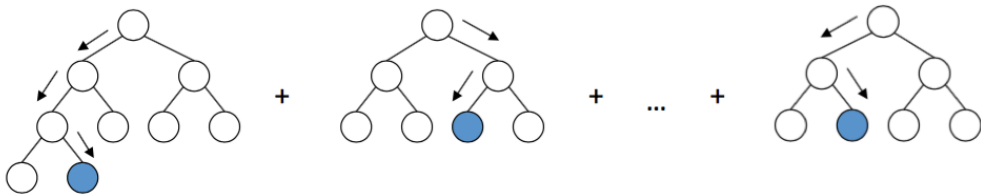- Each subsequent DT is trained to correct errors of previous DTs in the ensemble



**Figure:** BDT demonstration. Source: https://arogozhnikov.github.io/

# BDT: hyperparameters and benefits

Main tunable hyperparameters:

- **Max. depth**: The maximum depth of a tree (usually <12)
- **Learning rate**: This determines the impact of each tree on the final outcome (usually $\approx 0.1$)
- **Number of trees**: How many trees in ensemble

**Benefits**:

- Fast for training and prediction
- Easier to tune
- Minimalistic

# BDT: optimized set of features

BDT from XGBoost:

- Optimized **set of features** (sorted by *importance*):

  1. AccumCharge
  2. $R_{cht}$
  3. $z_{cc}$
  4. $pe_{std}$

  5. nPMTs
  6. $ht_{kurtosis}$
  7. $ht_{25\%-20\%}$
  8. $R_{cc}$

  9. $ht_{5\%-2\%}$
  10. $pe_{mean}$
  11. $J_{cht}$
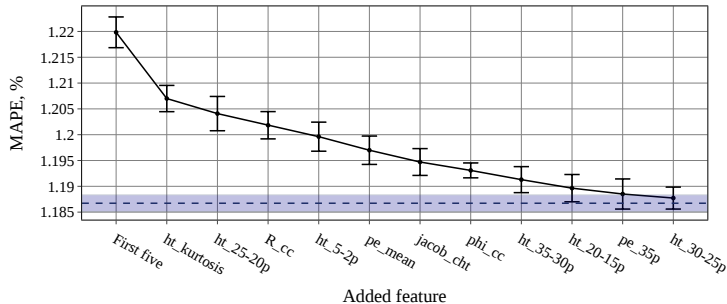  12. $\phi_{cc}$

  13. $ht_{35\%-30\%}$
  14. $ht_{20\%-15\%}$
  15. $pe_{35\%}$
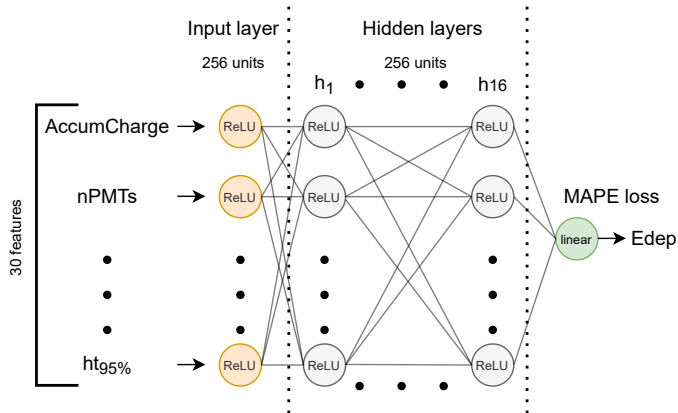  16. $ht_{30\%-25\%}$



- Optimized **hyperparameters** (using Grid Search):

  1. The maximum depth of the tree: 10
  2. Number of trees in the ensemble: $\simeq 500$
  3. Learning rate: 0.08

Fully-connected deep neural network (**FCDNN**):



- The search for hyperparameters was performed using *BayesianOptimizer*
- Training with *early stopping*
- Validation dataset: *400k events*
- *Selected features* provided the same performance as full set:

  1. AccumCharge
  2. nPMTs
  3. $R_{cc}$
  4. $R_{cht}$
  5. $\rho_{cc}$
  6. $\rho_{cht}$
  7. $pe_{mean}$
  8. $pe_{std}$
  9. $pe_{skew}$
  10. $pe_{kurtosis}$
  11. Percentiles of FHT distribution:
      $\{ht_{2\%}, ht_{5\%}, ht_{10\%}, ht_{15\%}, ..., ht_{90\%}, ht_{95\%}\}$
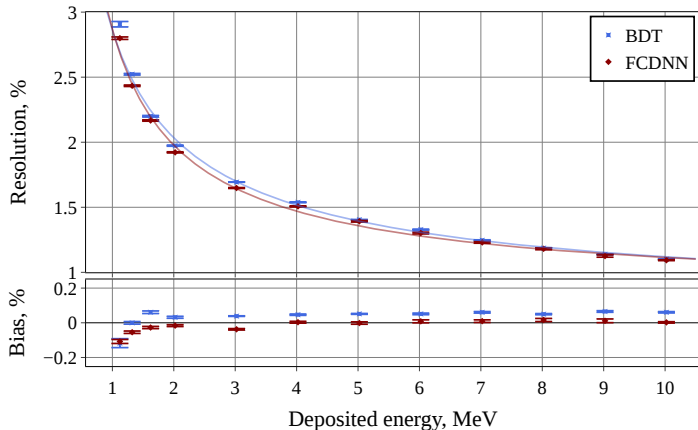
# Results

**Metrics**:

- Defined by a Gaussian fit of the $E_{\text{predicted}} - E_{\text{dep}}$ distributions
- *Resolution*: $\sigma / E_{\text{dep}}$, where $\sigma$ — standard deviation of the fit
- *Bias* $\mu / E_{\text{dep}}$, where $\mu$ — mean of the fit

**Parameterization**:

$$\frac{\sigma}{E_{\text{dep}}} = \sqrt{\left(\frac{a}{\sqrt{E_{\text{dep}}}}\right)^2 + b^2 + \left(\frac{c}{E_{\text{dep}}}\right)^2}$$

**Models' pred. time and memory usage**:

|                       | BDT | FCDNN |
|-----------------------|-----|-------|
| Pred. time, sec/100k  | 3.5 | 17    |
| Size, MB              | 50  | 12    |

# Summary

- **Energy reconstruction** using the information collected by PMTs

- *Aggregated* features approach

- The following ML models are used: **BDT, FCDNN**

- As a result <u>*achieved*</u>:
  1. High **quality** 3% @ 1 MeV, requared for physics goals of JUNO
  2. Great **computation speed**, thanks to a small set of aggregated features (in $10^4 - 10^5$ times faster than traditional methods)

# References and more details

Publications:

- **A. Gavrikov**, et al. arXiv: 2206.09040 (2022)
- **A. Gavrikov**, et al. EPJ Web Conf. 251 (2021), 03014
- Z. Qian, et al. NIM-A 1010 (2021), 165527