# gMCAP: a GPU-based Monte Carlo proton transport program for high-density tissues with precise nuclear reaction models

Xi-Yu Luo[1,2] · Liang Sun[3] · Zhen Wu[1,4] · Rui Qiu[1,2] · Shou-Ping Xu[5] · Hui Zhang[1,2] · Jun-Li Li[1,2]

## Abstract

GPU-based Monte Carlo (MC) simulations are highly valued for their potential to improve both the computational efficiency and accuracy of radiotherapy. However, in proton therapy, these methods often simplify human tissues as water for nuclear reactions, disregarding their true elemental composition and thereby potentially compromising calculation accuracy. Consequently, this study developed the program gMCAP (GPU-based proton MC Algorithm for Proton therapy), incorporating precise discrete interactions, and established a refined nuclear reaction model (REFINED) that considers the actual materials of the human body. Compared to the approximate water model (APPROX), the REFINED model demonstrated an improvement in calculation accuracy of 3%. In particular, in high-density tissue regions, the maximum dose deviation between the REFINED and APPROX models was up to 15%. In summary, the gMCAP program can efficiently simulate 1 million protons within 1 s while significantly enhancing dose calculation accuracy in high-density tissues, thus providing a more precise and efficient engine for proton radiotherapy dose calculations in clinical practice.

**Keywords** Monte Carlo simulation · Proton therapy · Dose calculation · GPU · Geant4

## 1 Introduction

Dose calculation remains a significant numerical challenge in radiation applications [1–3], particularly in particle radiotherapy, which is an interesting research area, including proton and ion radiotherapy [4–6], involving devices [7–10],

✉ Liang Sun
slhmz666@suda.edu.cn

Jun-Li Li
lijunli@mail.tsinghua.edu.cn

1 Department of Engineering Physics, Tsinghua University, Beijing 100084, China

2 Key Laboratory of Particle and Radiation Imaging, Tsinghua University, Ministry of Education, Beijing 100084, China

3 State Key Laboratory of Radiation Medicine and Protection, School of Radiation Medicine and Protection, Soochow University, Suzhou 215123, China

4 Nuctech Company Limited, Beijing 100084, China

5 National Cancer Center/ National Clinical Research Center for Cancer/Cancer Hospital, Chinese Academy of Medical Sciences and Peking Union Medical College, Beijing 100021, China

dose algorithm [11, 12], scanning modes [13], etc. Accurate dose algorithms are essential for maximizing the effectiveness of proton radiotherapy by leveraging the unique Bragg-peak characteristics of this approach [14]. Several dose algorithms have been proposed for proton radiotherapy, including the pencil beam algorithm [15–17], track-repeating methods [18, 19], the Boltzmann transport equation solver [20, 21], and full MC simulations [22, 23]. Although the pencil beam algorithm is widely used in clinical settings, there is an urgent need to improve its accuracy. The MC technique offers the highest accuracy, but it is hampered by prolonged computation times. Thus, the primary goal of proton dose calculation is to ensure both the precision and efficiency of the algorithm simultaneously.

In recent years, the integration of graphics processing units (GPUs) into photon dose algorithms has led to significant advancements, providing practical solutions through the implementation of various balancing strategies [24–28]. Similarly, GPU-based MC methods have shown promising results for proton radiation transport [29–33].

According to data from ICRU Report 46 [34], the weight ratio of hydrogen (H) atoms to the combined weight of nitrogen (N), carbon (C), and oxygen (O) atoms

in human tissues is believed to mirror the proportion of H atoms to O atoms in water. Furthermore, as outlined in ICRU Report 63 [35], these atoms are arranged in the same order when adjusting the nuclear reaction cross sections of C, N, and O atoms to their atomic masses. This alignment enables human tissues to behave similarly to water during nuclear reactions. Consequently, many proton programs assume that human tissues mimic the behavior of water, focusing primarily on proton interactions with H and O atoms in the nuclear reaction model [29, 32, 36].

In contrast, according to ICRP publication 70 [37], the skeleton comprises approximately 20% of the average total body weight, excluding adipose tissue, with 5% to 20% of the skeleton being composed of calcium (Ca) and a small amount of phosphorus (P). The nuclear interaction cross sections of these atoms are approximately 25% lower than that of oxygen (O) [38]. Consequently, substituting water for the elements present in the human body during dose calculations may lead to specific discrepancies. Additionally, atoms of elements such as sodium (Na), magnesium (Mg), phosphorus (P), sulfur (S), chloride (Cl), potassium (K), iron (Fe), and iodine (I) exhibit notable differences from O atoms. However, to date, only a few programs have developed extensive nuclear reaction models that consider interactions between protons and other constituent elements of human tissues [31].

Both types of programs have shown promising results and efficiency through their distinct strategies. Nevertheless, there is a dearth of quantitative research investigating the impact of the refinement level in nuclear reactions on dose calculations for human tissues. Given that increased refinement in nuclear reactions increases algorithmic effort and redundancy, comprehensive research into this topic is warranted.

This study aimed to develop an advanced proton transport program named gMCAP that utilizes a GPU. The code incorporates precise discrete interactions and refined nuclear reaction models. The initial implementation of gMCAP refined electromagnetic interactions using approaches such as the energy loss method and multiple scattering model. Subsequently, two nuclear reaction models were developed: simplified (APPROX) and refined (REFINED). The APPROX model considers proton interactions with water of varying densities, encompassing p-H elastic nuclear reactions, p-O elastic nuclear reactions, and p-O inelastic nuclear reactions. In contrast, the REFINED model accounts for interactions between protons and specific elements present in human tissues, including H, C, N, O, Na, Mg, P, S, Cl, K, Ca, Fe, and I atoms. Finally, the entire program was ported to the CUDA platform to leverage parallel architecture, and the accuracy and computational efficiency of the code were assessed.

## 2 Materials and methods

In the gMCAP code, both APPROX and REFINED models utilize the same electromagnetic interaction algorithm. This article focuses on modeling and correction of electromagnetic interactions in Subsections A-C. Subsection D introduces the nuclear reaction algorithm employed by the APPROX model, specifically addressing the nuclear interactions between protons and H and O atoms. In Subsection E, the nuclear reaction algorithm of the REFINED model is presented, encompassing elastic nuclear reactions and inelastic nuclear reaction processes occurring between protons and each element present in the human body.

### 2.1 Geometry and proton transport algorithm

Serial CT images are routinely used for the planning and optimization of clinical treatment. First, standard calibration methods such as the Hounsfield unit (HU) look-up table are utilized to transform CT images into a three-dimensional density matrix $\rho$, as illustrated in Fig. 1a. Subsequently, regions of interest (ROI) are delineated using either threshold segmentation or artificial intelligence-based segmentation algorithms [39] to accurately represent tissues and organs with distinct material compositions. Furthermore, to determine the energy deposition at a specific step length $\Delta s$, we utilized the stopping power ratio (SPR) to rescale it to a step length in water, referred to as $\Delta s_w$, using Eq. (1). The SPR is related to the proton energy and material density. It can be determined by using a fitting formula [22], stoichiometric calibration [40], or artificial intelligence [41]. In this case, the first method was utilized, and a portion of the conversion curve (160 MeV) is shown in Fig. 1b. The geometric structure and processes employed enable the transport of protons within voxels.
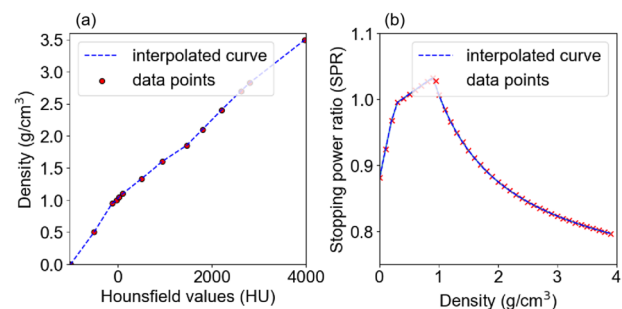


**Fig. 1** (Color online) Material look-up table used in gMCAP: **a** Hounsfield value look-up table to determine density from CT images (for different CT machines, the curves are different), and **b** stopping power relative to water of different densities of material for 160 MeV protons

$$\Delta s_{\mathrm{w}} = SPR\frac{\rho}{\rho_{\mathrm{w}}}\Delta s \tag{1}$$

The transport of protons in the medium is influenced by both elastic and inelastic Coulomb collisions as well as elastic and inelastic nuclear interactions. Among these, numerous elastic Coulomb collisions occur, with almost $1\times10^6$ instances per centimeter of the path [22]. Because of the impracticality of simulating these interactions individually, a Class II condensed history algorithm was employed for proton transport. A production threshold $T_{\mathrm{e}}^{\min}$ was established, which indicates the minimum energy required for ionization to occur and generate $\delta$ electrons. If the energy is below $T_{\mathrm{e}}^{\min}$, continuous energy loss occurs instead. Thus, reactions between protons and materials can be divided into two categories: continuous interactions (see Sect. 2.2) and discrete interactions, including ionization (see Sect. 2.3), elastic nuclear reactions (see Sect. 2.5.1), and inelastic nuclear reactions (see Sect. 2.5.2).

Protons are transported stepwise in the voxel geometry. One step $\Delta s$ is defined by the minimum of three step sizes: (1) the distance to the closest voxel boundaries; (2) the user-defined maximum step size, which is set to 1 mm, or the residual range corresponding to 20% of the initial energy; and (3) the mean free path of each interaction, including ionization, elastic nuclear interaction, and inelastic interaction. The mean free path $\lambda$ is defined by Eq. (2):

$$\lambda = -\log(\eta)\frac{1}{\Sigma}. \tag{2}$$

Here, $\eta$ is a random number between 0 and 1, and $\Sigma$ is the cross section of each discrete interaction.

## 2.2 Continuous interactions

### 2.2.1 Energy loss algorithm

According to Zhang et al. [42], the energy loss ($\Delta E$) of a geometric step $\Delta s$ in a material is equal to the energy deposited in the step $\Delta s_{\mathrm{w}}$ in water. The proton's energy loss fluctuates around its mean energy loss $\overline{\Delta E}$, and $\overline{\Delta E}$ for a step is generally determined by integrating the restricted stopping power $L_{\mathrm{w}}$, as shown in Eq. (3). $L_{\mathrm{w}}$, which is consistent with $T_{\mathrm{e}}^{\min}$, is calculated using Geant4 [43] and is dependent on the proton energy $T_{\mathrm{p}}$. Nevertheless, the process of obtaining the integral solution typically requires a significant amount of time, which is contingent upon the magnitude of the differential step size d$x$.

$$\Delta s_{\mathrm{w}} = -\int_{T_{\mathrm{p}}}^{T_{\mathrm{p}}-\overline{\Delta E}} \frac{\mathrm{d}T_{\mathrm{p}}'}{L_{\mathrm{w}}\left(T_{\mathrm{p}}', T_{\mathrm{e}}^{\min}\right)} \tag{3}$$

Therefore, a proton Energy Range LookUp Table (ERLUT) is proposed to solve this issue. Initially, the $L_{\mathrm{w}}$ values of protons with various energies in water are determined. Subsequently, the $L_{\mathrm{w}}$ values of protons with varying energies are integrated to obtain their range in water. Finally, an ERLUT derived from the precise integration is acquired, as shown in Fig. 2a.

For a proton with initial energy $T_{\mathrm{p}}$, the energy loss $\overline{\Delta E}$ in step $\Delta s$ can be obtained by interpolating the ERLUT twice. First, the remaining range corresponding to $T_{\mathrm{p}}$ is obtained via forward interpolation. The table is then inversely interpolated to obtain $T_{\mathrm{p}}^{\mathrm{end}}$ at the end of the step. The difference between the remaining ranges corresponding to $T_{\mathrm{p}}$ and $T_{\mathrm{p}}^{\mathrm{end}}$ is the step size $\Delta s$. Using this approach can not only save time but also ensure precision. Figure 2b, c shows a comparison between the continuous slowing down approximation (CSDA) range and the calculation time for 200 MeV protons. The comparisons
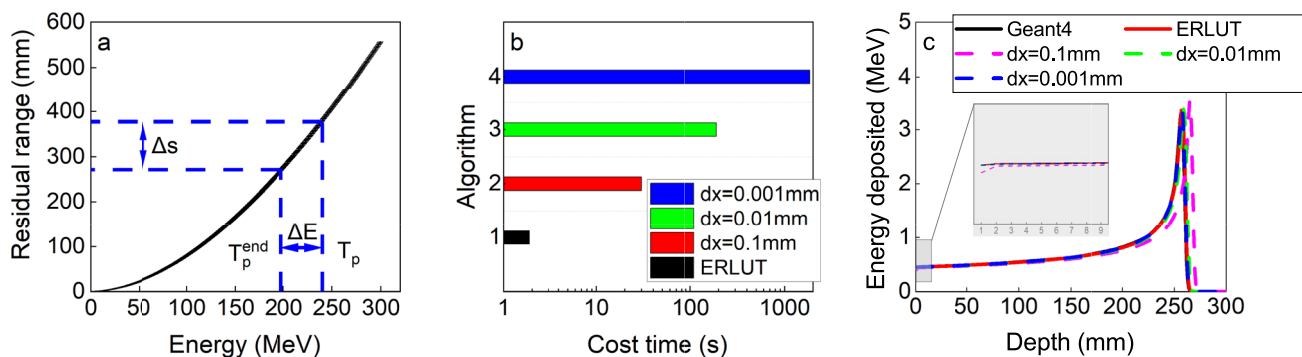


**Fig. 2** (Color online) **a** The remaining range of protons with various energies (the ERLUT), **b** the cost time of protons with 200 MeV under different integration precisions, and **c** the CSDA range of protons with 200 MeV under different integration precisions

were performed under different integration precisions ($dx$) utilizing the ERLUT.

The energy straggling model was implemented to rectify the exact energy deposition of the protons within a single step. The energy deposition exhibits a Gaussian distribution with $\overline{\Delta E}$ and a variation of $\sigma$, which is defined in Eq. 4.

$$\sigma^2 = 2\pi r_e^2 m_e n_e \Delta s \frac{\min\left(T_e^{\min}, T_e^{\max}\right)}{\beta^2}\left(1 - \frac{\beta^2}{2}\right) \tag{4}$$

### 2.2.2 Multiple scattering model

Protons undergo elastic Coulomb collisions when they traverse matter, resulting in numerous small-angle deflections. These deflections can be simulated using a multiple scattering model. This study employs the multiple scattering theory proposed by Rossi and Greisen [44], which states that the deflection angle on the projection plane containing the initial track obeys a Gaussian distribution with $\sigma$ as shown in Eqs. (5)-(6):

$$\sigma = \frac{1}{\sqrt{2}} \times \sqrt{\frac{1}{2}E_s^2 t/p^2\beta^2}, \tag{5}$$

$$t \equiv \frac{\Delta s}{X_0(\rho)}. \tag{6}$$

Here, $p$ is the momentum of a proton, $\beta$ is its velocity, $\Delta s$ is the step length, and $X_0(\rho)$ is the radiation length of the material. According to Rossi and Greisen, the mean square angle of scattering is independent of the atomic number, and $E_s$ is a constant parameter of the multiple scattering model with the dimension of energy, given by Eq. (7):

$$E_s = \mu_e(4\pi 137)^{\frac{1}{2}} = 21\,\text{MeV}. \tag{7}$$

Here, $\mu_e$ is the mass of the electron. Note that the value of $E_s$ here differs from that of Fippel [22] because in Fippel's model, electrons are considered to slow down continuously.

Within the gMCAP program of the registration-only electromagnetic interaction model, comparisons were made between the central axis depth dose (CADD) curve and the lateral dose distribution (LDD) at the peak position across various $E_s$ values. The experimental results, illustrated in Fig. 3, were compared with Geant4 simulations to identify the most appropriate Es value. After conducting the simulations, the energy parameter was refined by setting $E_s$ to 23.5 MeV. This adjustment was implemented to maximize the precision of the electromagnetic process.
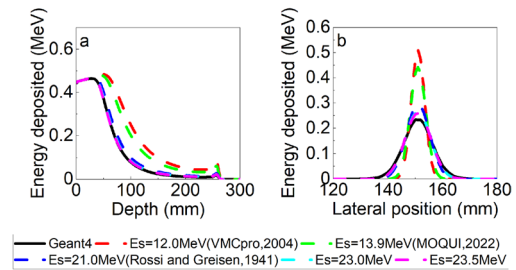
**Fig. 3** (Color online) **a** Central axis depth dose curves and **b** lateral dose curves in the peak for different $E_s$ values for a cube of water irradiated by 200 MeV protons
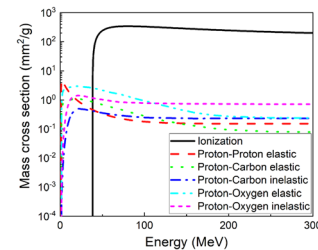


**Fig. 4** (Color online) Cross-sectional data of discrete interactions in G4_SKIN_ICRP taken from Geant4

### 2.3 Ionization

The ionization interactions between protons and matter are crucial for dose estimation. The cross section for electromagnetic interactions is approximately two orders of magnitude larger than that for nuclear reactions. Take G4_SKIN_ICRP as an example, which is composed of H (10%), C (20.4%), N (4.2%), O (64.5%), Na (0.2%), P (0.1%), S (0.2%), Cl (0.3%), K (0.1%). Figure 4 shows the cross-sectional data obtained from Geant4 scaled by mass density.

According to the two-body collision energy transfer relationship, the maximum energy $T_e^{\max}$ transferred by a proton to a free electron is calculated using Eq. 8.

$$T_e^{\max} = \frac{2m_e\beta^2\gamma^2}{1 + 2\gamma m_e/m_p + (m_e/m_p)^2} \tag{8}$$

Here, $m_e$ and $m_p$ are the electron and proton rest masses, respectively; the relativistic parameter $\gamma = \frac{T_p + m_p}{m_p}$; and $\beta^2 = 1 - \frac{1}{\gamma^2}$. The differential macroscopic cross section for ionization to produce a $\delta$-electron is extracted from Geant4 using the class G4hIonisation with $T_e^{\min} = 81.5\,\text{keV}$. This corresponds to a range of approximately 0.1 mm in water. The G4hIonisation class utilizes the following models:

- G4BetheBlochModel, valid for $T_p > 2$ MeV

- G4BraggModel, valid for $T_p < 2$ MeV

If an ionization event occurs during transport, a $\delta$-electron is produced. The energy of the $\delta$-electron is determined using a random sampling method with a probability distribution function [22]. Figure 5 shows the electron energy distribution produced by a 200 MeV proton.

The scattering angular deflection of the proton during this process is ignored, as the mass of the proton is much larger than that of an electron, and its direction hardly changes during a collision. The $\delta$-electrons deposit energy locally and cease to be transported.

## 2.4 Simplified nuclear reaction: APPROX model

The simplified APPROX model focuses on nuclear interactions between protons and H and O atoms, including elastic and inelastic reactions. In the APPROX model, soft tissue is presumed to undergo nuclear reactions similar to those in water. Assuming this, all tissues consist of H and O with a mass ratio of 1:8. Different tissues and organs exhibit variations in density.

The cross sections of the three nuclear reactions were calculated using Geant4 and tabulated. The model proposed by Fippel and Soukup [22] was used to sample the energy and scattering angle of the nucleus. For the secondary particles, where secondary protons are stored in the stack for further transportation, other charged particles release energy locally, and long-range particles, such as photons and neutrons, deposit energy proportionally.

## 2.5 Refined nuclear reaction: REFINED model

Because human tissues contain atoms other than H and O, it is necessary to consider the nuclear reactions of these atoms. Particularly for middle-$Z$ and high-$Z$ atoms, such as P, Ca, Fe, and I, the normalized nuclear reaction cross section with
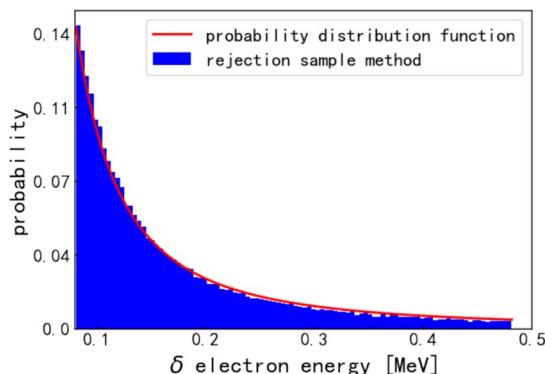


**Fig. 5** (Color online) The $\delta$-electron energy distribution produced by a 200 MeV proton

protons is significantly lower than that of O. Thus, employing the APPROX model to estimate the radiation dose to the bone may result in inaccuracies.

Consequently, a more advanced nuclear reaction model, REFINED, was developed in this study. It encompasses the nuclear interactions between protons and elements including H, C, N, O, Na, Mg, P, S, Cl, K, Ca, Fe, and I. Similarly, the cross sections of protons with these atoms were obtained using Geant4. Next, the program was initialized to compute the cross sections for various materials and each individual nuclear reaction. At the post-step point, the discrete interaction type was sampled and specific nuclear reaction procedures were executed.

### 2.5.1 Elastic nuclear interactions

An elastic nuclear reaction is a two-body collision event that follows the laws of conservation of energy and momentum. The deflection angle of the recoil nuclei is sampled based on the differential section and calculated using Ranft's empirical formula [45], as shown in Eq. 9:

$$
\begin{cases}
\dfrac{d\sigma_{el}}{d\Omega} \approx A^{1.63} \exp(14.5A^{0.66}t) + 1.4A^{0.33}\exp(10t) \\
\quad A \leq 62, \\
\dfrac{d\sigma_{el}}{d\Omega} \approx A^{1.33}\exp(60A^{0.33}t) + 0.4A^{0.40}\exp(10t) \\
\quad A > 62.
\end{cases}
\tag{9}
$$

Among them, the invariant momentum transfer $t \equiv -2p^2(1 - \cos\theta)$, the unit is $(\text{GeV}/c)^2$, and $p$ and $\theta$ are the momentum and deflection angle of the recoil nuclei in the center-of-mass system, respectively. Normalization of the differential cross-sectional formula is not important for the random selection of the scattering angle. A rejection method [36] is used to sample $t$, after which the angle $\theta$ is obtained. Subsequently, the Lorentz boost [28], which is used in gMCAP, is used to determine the angle and momentum in the laboratory system.

Once the deflection angle $\theta$ has been calculated, the energies of the incident proton and recoil nuclei can be acquired, and the direction of the incident proton can be updated. The energy of the heavy recoil nucleus is assumed to be deposited locally, and the updated protons continue to be transported.

### 2.5.2 Inelastic nuclear interactions

Inelastic nuclear interactions cause the proton to evaporate and transfer its energy to secondary particles. Similarly, distinct reaction channels are sampled depending on the cross section of the inelastic nuclear reactions between the protons

and each atomic nucleus. This study rectifies the distribution of the scattering angles for secondary protons, considering various target nuclei. According to Qin et al. [36], the sample expression for the angular cosine value $\cos\theta$ of the secondary proton is given as Eq. 10:

$$\cos\theta = \ln[(e^{\zeta} - e^{-\zeta})\eta + e^{-\zeta}]\frac{1}{\zeta}. \tag{10}$$

Here, $\eta$ is a random number between 0 and 1. Maneval et al. [31] described the parameter $\zeta$ as shown in Eq. 11:

$$\zeta = \frac{T_p}{T^{\alpha}}, \text{ and } \alpha = \mu[1 + \tau \ln T]. \tag{11}$$

Here, $T$ denotes the energy of the secondary proton in MeV. Equation 12 expresses the parameter $\mu$, which represents the correction for the target nucleus and is dependent on the mass number $A$ and $T_p$.

$$\begin{cases} \mu = 1.8 \times 10^{-2}(A - 12) - 10^{-3}T_p + 3.69, \\ T_p < 100\,\text{MeV}, \\ \\ \mu = 0.55 \times 10^{-2}(A - 12) + 10^{-3}T_p + 3, \\ 100\,\text{MeV} \le T_p \le 200\,\text{MeV}, \\ \\ \mu = 1.1 \times 10^{-2}(A - 12) + 10^{-3}T_p + 3, \\ T_p > 200\,\text{MeV}. \end{cases} \tag{12}$$

Another parameter, $\tau$, is defined as a quantity associated with $T_p$, given by Eq. 13.

$$\begin{cases} \tau = 10^{-3}T_p - 0.28, & T_p < 100\,\text{MeV}, \\ \tau = 0.2 \times 10^{-3}T_p - 0.203, & T_p \ge 100\,\text{MeV}. \end{cases} \tag{13}$$

In addition, the azimuth angle $\phi$ is uniformly sampled between $[0, 2\pi]$.
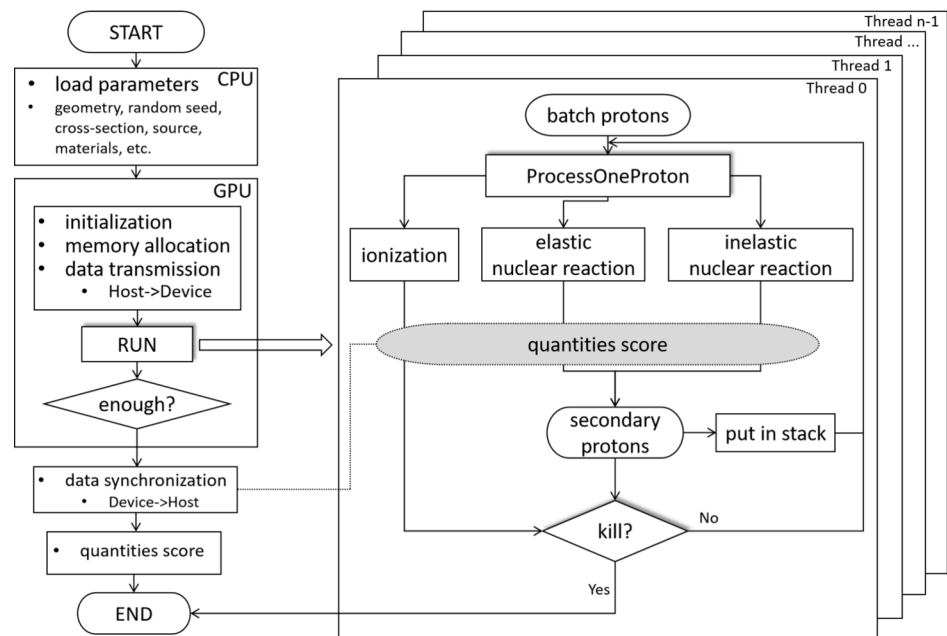
The approach proposed by Fippel and Soukup [22] allows for the classification of secondary particles into three distinct categories: secondary protons, large-range particles (such as photons and neutrons), and short-range particles (such as alpha particles and heavier nucleus fragments). The energy of the secondary particles is sampled from a uniform distribution between the minimum energy (3 MeV here) and the remaining energy of the system. Subsequently, the new remaining energy is assigned by subtracting the energy of the secondary particles and additional binding energy. The energy of the short-range particles is deposited locally, whereas that of the long-range particles is released proportionally.

## 2.6 GPU implementation

GPU parallel computing was implemented using CUDA v11.5, which supports programming in C/C++. When the program starts, the CPU initializes the cross-sectional data of the material, source information, geometric configuration, and random number seed. All these data are transferred from the CPU memory to the GPU global memory. After initialization, the primary and secondary protons are transported in batches, and each primary proton is transported using a GPU thread. Atomic addition is used to write the results to global memory to obtain the average dose of each voxel, thus avoiding memory conflicts and saving execution time.

The code employs a batch-based approach to compute statistical fluctuations, also known as statistical uncertainty. Typically, a single simulation consists of 10 batches. The



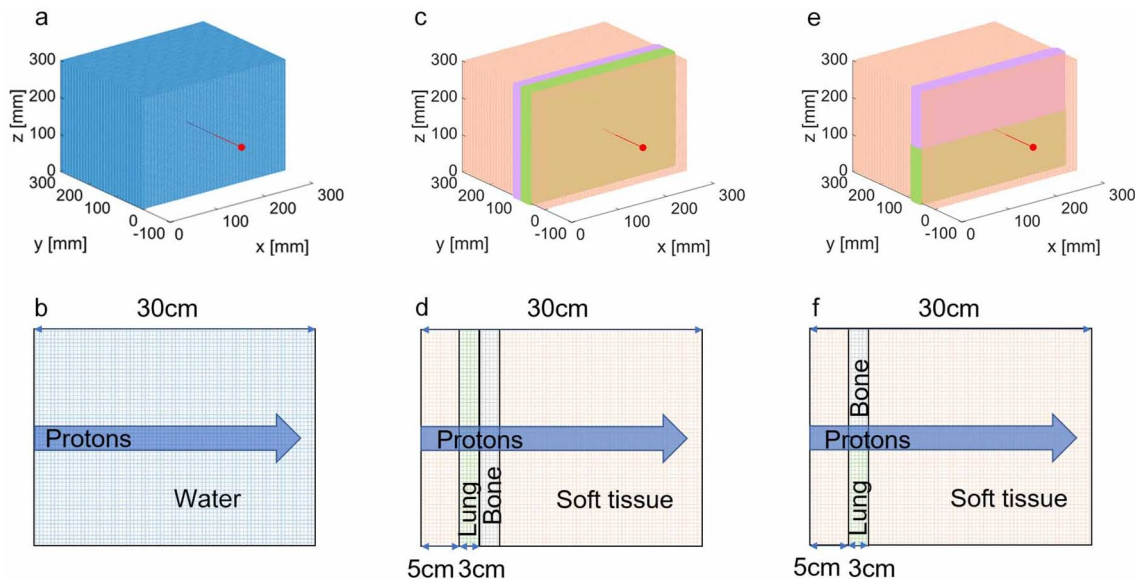**Fig. 6** (Color online) The CPU-GPU coupling architecture diagram of the entire gMCAP

**Fig. 7** (Color online) The homogeneous and heterogeneous phantoms and source: **a**, **b** show the schematic diagram of water, and **c–f** show the schematic diagram of the mixed model including soft tissue, lung tissue, and bone

calculation process is shown in Eq. 14, where $v$ is the voxel, $N$ is the number of batches, $X_i(v)$ is the physical quantity (dose or deposited energy) recorded by voxel $v$ in the $i$-th batch, and $\bar{X}(v)$ is the average physical quantity within the voxel.

$$s(v) = \sqrt{\frac{\sum_{i=1}^{N}(X_i(v) - \bar{X}(v))^2}{N(N-1)}} \tag{14}$$

Before the program exits, it copies the results from the GPU to the CPU and outputs the data. The overall program architecture is illustrated in Fig. 6.

## 2.7 Validation and comparison

### 2.7.1 Geant4 setup

Geant4 simulations were conducted using version 10.06 as a benchmark to compare dose distributions. The electromagnetic model utilized was G4EmStandardPhysics_option4, whereas G4HadronPhysicsQGSP_BERT_HP was employed as the hadron physics model along with G4HadronElasticPhysicsHP. In addition, the physical models G4StoppingPhysics, G4IonBinaryCascadePhysics, and G4DecayPhysics were included, as they are commonly used in proton radiation therapy [32]. The electron production cut was set to 0.1 mm, corresponding to an energy cutoff of 81.5 keV in gMCAP. The proton cutoff step size was set to 0.4 mm to balance accuracy and computation time. The Geant4 simulation was executed using an Intel Core i7-10875 H CPU operating at a clock speed of 2.30 GHz.

### 2.7.2 Homogeneous and heterogeneous phantoms

An extremely small proton beam was used to evaluate the stability of the gMCAP. The code was tested using three different phantoms: a uniform water phantom (WaterPhan), two standard heterogeneous phantoms (PhantomI and PhantomII) consisting of soft tissue, bone, and lung tissue. The voxel resolution for both the homogeneous and heterogeneous models was set to 1 mm × 1 mm × 1 mm, and the phantom size was 300 mm × 300 mm × 300 mm. In this experiment, proton beams with energies of 100 MeV, 150 MeV, and 200 MeV were employed to examine the water phantom. Additionally, a monoenergetic proton beam of 160 MeV was used to assess the heterogeneous phantom. Figure 7 illustrates a schematic of the source and geometry, as shown by Maneval et al. [31].
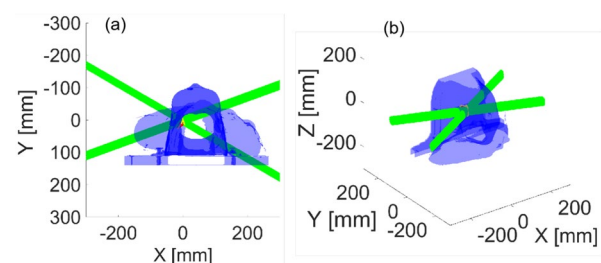


**Fig. 8** (Color online) Schematic diagram of the treatment plan in the laboratory coordinate system: **a** top view, **b** 3D view
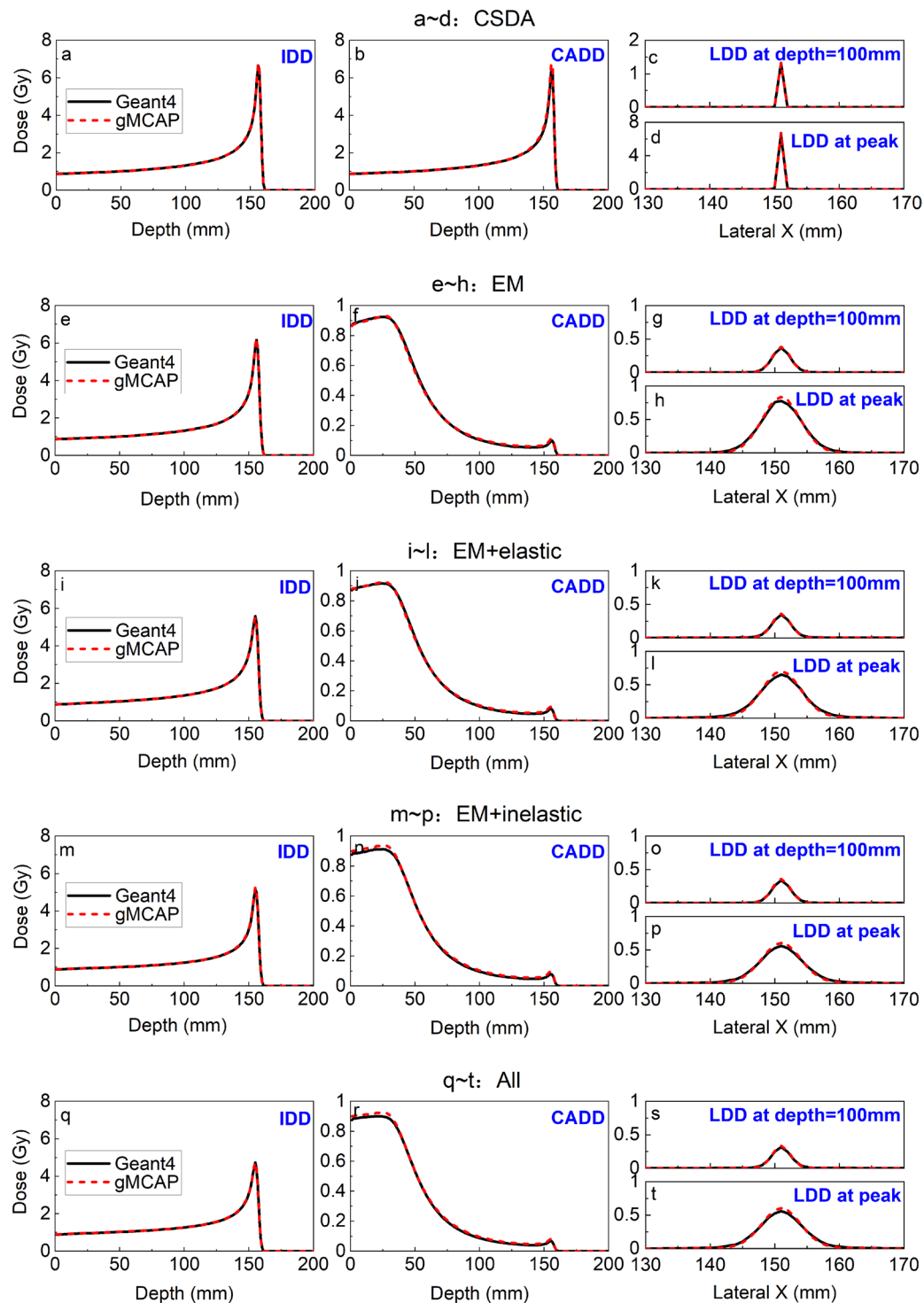
**Fig. 9** (Color online) Discrete interactions in water. The first column is the IDD curve, the second column is the CADD, and the third column is the LDD at depth = 100 mm and LDD at the peak. In the sub-graph, **a–d** are the results of CSDA, **e–h** are the results of EM, **i–l** are the results of EM and elastic interactions, **m–p** are the results of EM and Inelastic interactions, and **q–t** are the results of all interactions

### 2.7.3 Clinical cases based on patient treatment plans

The suitability of the clinical case was evaluated using a dataset comprising head and neck CT data. The CT data were voxelized at a resolution of 0.875 mm × 0.875 mm × 3.000 mm, resulting in a total size
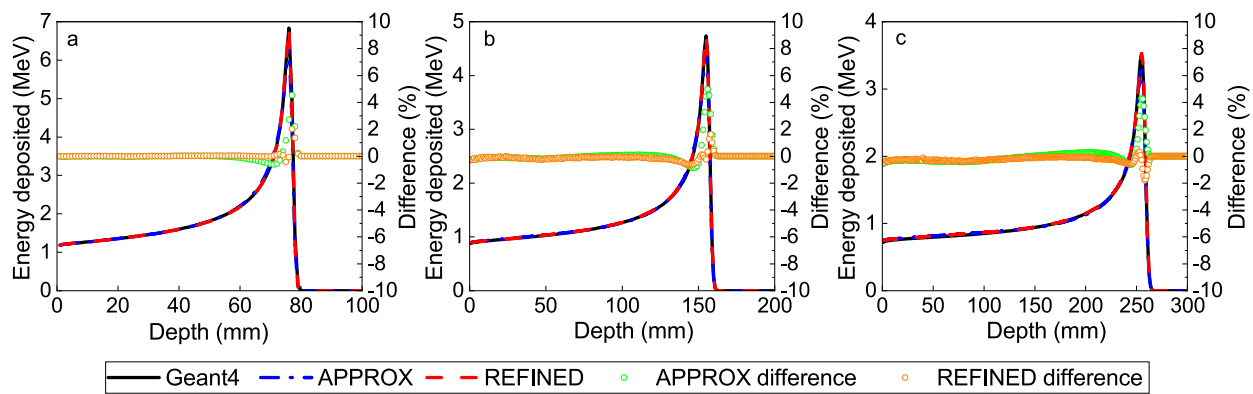
**Fig. 10** (Color online) Integrated depth dose in water with primary proton energy of **a** 100 MeV, **b** 150 MeV, and **c** 200 MeV

of $512 \times 512 \times 79$ voxels. Initially, a conversion process was applied to transform the HU values of the CT data into density values. Subsequently, the tissues were assigned. For the dose computation, the proton bixel width was set to 1 mm using the open-source treatment planning system matRad [46]. Information on 9651 rays from two irradiation fields was obtained. Figure 8 presents a schematic of the treatment plan in the laboratory coordinate system.

Moreover, in order to evaluate the impact of the two nuclear reaction models (APPROX and REFINED), proton beams with energies of 120 MeV, 140 MeV, and 180 MeV were utilized to irradiate the mandibular and tooth regions. These anatomical structures comprise tissues and organs that contain elements such as P and Ca. The dose distribution was recorded, and the relative deviations were analyzed using the two models.

# 3 Results

## 3.1 Validation for discrete interactions

The discrete interaction was validated by comparing the dose distribution of a 150 MeV proton beam interacting with water using different physical models. In this study, the refined physical model, REFINED, of gMCAP was employed, whereas the well-studied simplified model, APPROX, was omitted because it has been extensively examined by previous researchers. Various physical models were tested separately, including CSDA, Electromagnetic Interactions (EM, which includes CSDA, Multiple Scattering, and Energy Straggling), Elastic Nuclear Interactions (EM + Elastic), Inelastic Nuclear Interactions (EM + Inelastic), and the complete set of physical models (All).

To verify the accuracy of the energy loss and angular deflection processes in each physical model, the integral depth dose (IDD), central axis depth dose (CADD), and

lateral depth dose (LDD) at depth = 100 mm with LDD at the peak were examined and compared with the corresponding Geant4 results. The results, depicted in Fig. 9, demonstrate a strong agreement between gMCAP and Geant4, confirming the accuracy of the program's physical model. Notably, the LDD at the peak in gMCAP is slightly higher than that in Geant4. This discrepancy could be attributed to differences in the multiple scattering model being utilized. Despite the corrections made to this model in the present study, some deviations from Geant4 still exist. This discrepancy may be attributed to the exclusion of electron transport. However, it is important to emphasize that the level of inaccuracy is within an acceptable level. Further research on the transport model of secondary electrons will be conducted in the future.

## 3.2 Quantitative analysis

A quantitative analysis was conducted to compare the dose distributions of protons at different energy levels. This analysis involved the use of the APPROX and REFINED nuclear reaction models. The IDD curve was examined, as shown in Fig. 10. The results indicated that the disparity between the APPROX model and Geant4 was less than 5%. However, the REFINED model further reduced the overall deviation to 2%, with only a 1% deviation observed in the region beyond the distal fall-off area.

The analysis of the lateral distribution is presented in Fig. 11. The results demonstrate that the relative deviation between Geant4 and 90% of the data fell within 1%, indicating that the REFINED model exhibited better performance than the APPROX model.

In addition, the heterogeneous model was used to evaluate the IDD curve and lateral dose profiles; the results are illustrated in Fig. 12. The results demonstrate good agreement between the dose distributions obtained from gMCAP and Geant4. Specifically, the relative deviation in the lateral dose
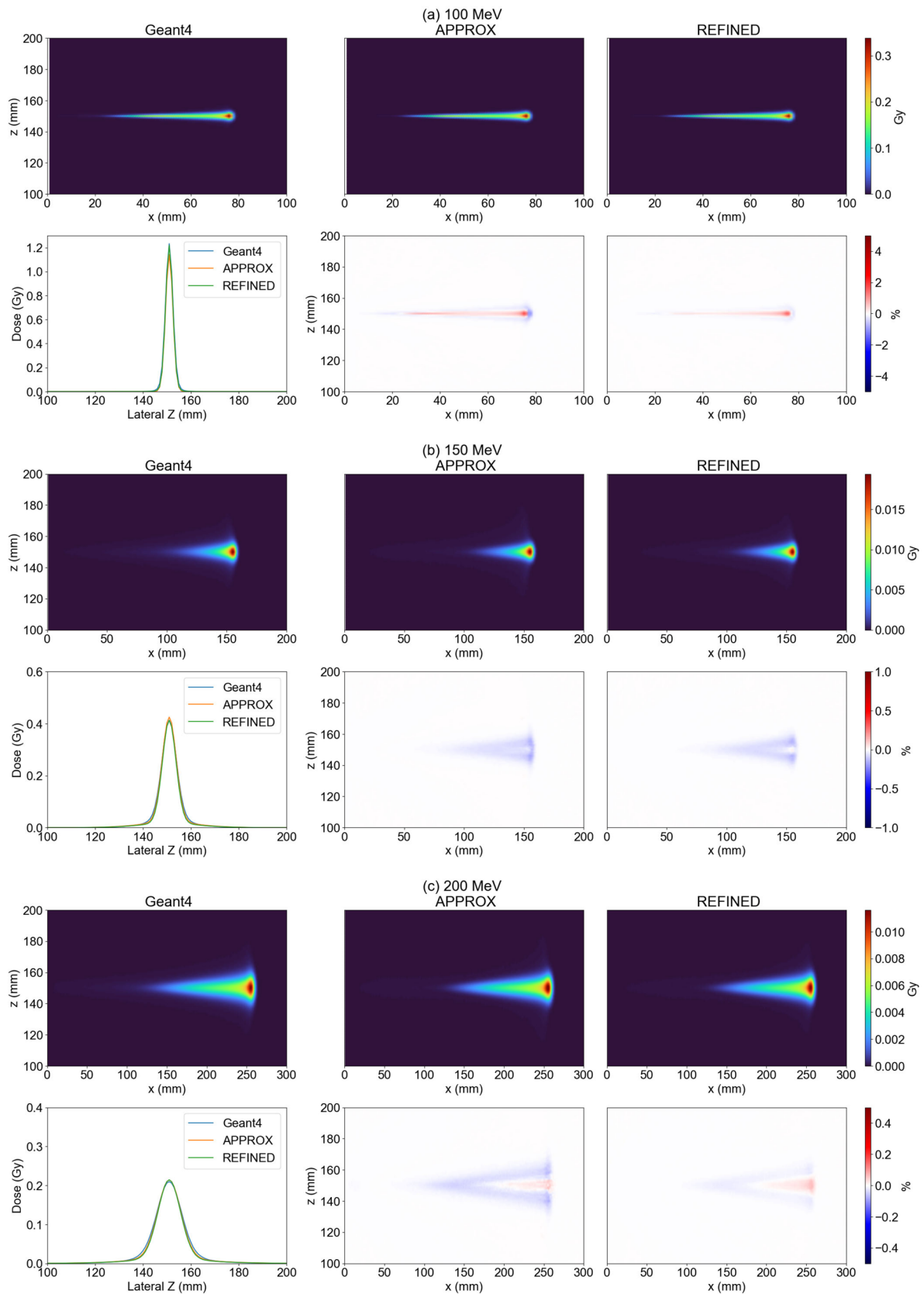
◀**Fig. 11** (Color online) Lateral dose distribution in water with primary proton energy of **a** 100 MeV, **b** 150 MeV, and **c** 200 MeV. The first row of each subfigure shows the dose distributions of Geant4, APPROX, and REFINED, respectively. The second row compares the lateral distributions at the peak, the relative error distributions between APPROX and Geant4, and the relative error distributions between REFINED and Geant4. The relative error is calculated as $\text{error}_i = \left(\text{dose}_i^{\text{APPROX|REFINED}} - \text{dose}_i^{\text{Geant4}}\right) \Big/ \max\left(\text{dose}^{\text{Geant4}}\right)$

profile was found to be within 3%. Notably, the REFINED model exhibited a reduced overall deviation compared to Geant4, in contrast to the APPROX model.

The gamma passing rate, defined in [47], is a widely employed metric in radiotherapy for assessing the similarity of dose distribution maps. It considers voxels with doses greater than 10% of the maximal dose. The findings, presented in Table 1, indicate a high level of concurrence with the Geant4. The gamma passing rate of 2 mm/2% exceeds 99%. Moreover, the REFINED model outperformed the APPROX model in terms of results.

### 3.3 Dose results of clinical cases

The dose outcomes for treatment plans targeting the head and neck region were compared. The dose distribution for the middle slice is shown in Fig. 13. The overall dose results closely match those of Geant4, with the relative deviation

between gMCAP and Geant4 within 4%. Moreover, the gamma passing rates for both the APPROX and REFINED models were nearly identical at approximately 99%.

Furthermore, the dose outcomes along a line intersecting both beams and parallel to the *y*-axis were analyzed. The corresponding gamma values are also presented in the provided figure. This analysis provides additional insight into the agreement between gMCAP and Geant4 in terms of dose accuracy and conformity.

### 3.4 Efficiency gain

Compared to Geant4 running on a CPU for over 5 h, gMCAP offers a significant reduction in runtime, completing the simulations within seconds. This translates to an efficiency gain of more than 1800 times in comparison with Geant4. The execution times for the two versions of gMCAP using the APPROX and REFINED models are summarized in Fig. 14. Using an NVIDIA GeForce RTX 4090 GPU, gMCAP achieved the transportation of one million particles in less than 1 s. Note that the REFINED model requires approximately 1.15 times more time than the APPROX model, with variations depending on the geometry and source type used in the simulations. These results highlight the efficiency and speed advantages of gMCAP for particle transport simulations.

**Fig. 12** (Color online) Integrated depth dose and lateral dose profile in heterogeneous model with primary proton energy of 160 MeV
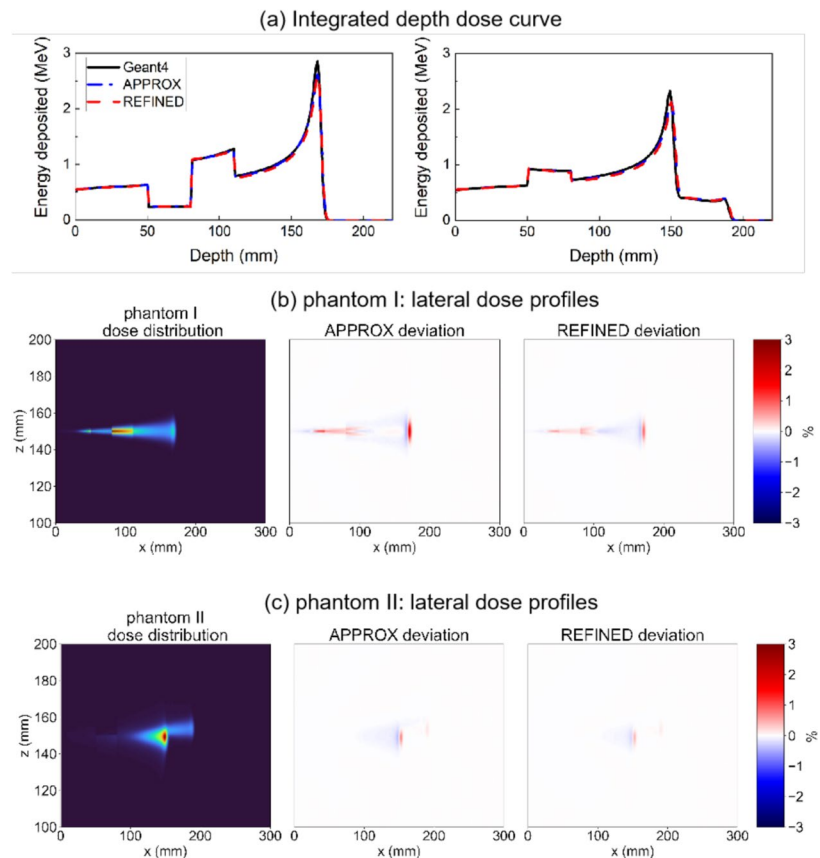
**Table 1** Gamma passing rates for the homogeneous and heterogeneous phantoms

| Geometry | Criteria | APPROX (%) | REFINED (%) |
|---|---|---|---|
| WaterPhan 200 MeV | 1 mm, 1% | 98.39 | 99.26 |
| | 2 mm, 2% | 98.90 | 100 |
| | 3 mm, 3% | 99.94 | 100 |
| PhantomI 70 MeV | 1 mm, 1% | 97.14 | 97.75 |
| | 2 mm, 2% | 98.97 | 99.31 |
| | 3 mm, 3% | 99.82 | 99.97 |
| PhantomI 160 MeV | 1 mm, 1% | 97.63 | 98.12 |
| | 2 mm, 2% | 99.46 | 99.90 |
| | 3 mm, 3% | 100 | 100 |
| PhantomII 70 MeV | 1 mm, 1% | 95.15 | 95.28 |
| | 2 mm, 2% | 98.16 | 99.53 |
| | 3 mm, 3% | 99.27 | 99.46 |
| PhantomII 160 MeV | 1 mm, 1% | 95.33 | 95.73 |
| | 2 mm, 2% | 98.88 | 99.87 |
| | 3 mm, 3% | 99.86 | 99.93 |



**Fig. 13** (Color online) Dose distribution comparison of the head and neck treatment plan. The first row shows the dose distributions of Geant4, APPROX, and REFINED, respectively. The second row shows the integrated lateral dose distribution, the relative error distribution between APPROX and Geant4, and the relative error distribution between REFINED and Geant4. The third row shows the gamma index distribution of APPROX and REFINED compared to Geant4, respectively

## 3.5 The impact of various nuclear models

Figure 15 shows the dose distribution and relative deviation of the two nuclear reaction models. Deviation primarily occurs in the teeth and downstream regions. The proton energies of 120 MeV, 140 MeV, and 180 MeV result in a maximum variance of approximately 15%. The primary explanation for this deviation is the ability of nuclear events to substantially alter the angular dispersion of secondary particles, primarily secondary protons. The tooth area comprises approximately 29% Ca and 14% P. The angular distribution of secondary protons generated through nuclear reactions differs from that produced by p-O nuclear reactions, leading to distinct energy and momentum distributions of the secondary protons downstream. The REFINED model incorporates the nuclear reactions of $Z$ elements such as P and Ca, which rectify the angular distribution of secondary particles, leading to distinct downstream dose distributions.

## 4 Conclusion and discussion

We have developed a GPU-accelerated program, gMCAP, which accurately calculates Monte Carlo (MC) proton dose distributions with precise discrete interactions and refined nuclear reaction models. gMCAP offers two nuclear interaction models: APPROX and REFINED. The APPROX model considers only H and O atoms, whereas the REFINED model considers all constituent elements of human tissues.

gMCAP with the REFINED nuclear model provides more accurate proton transport compared to that of certain
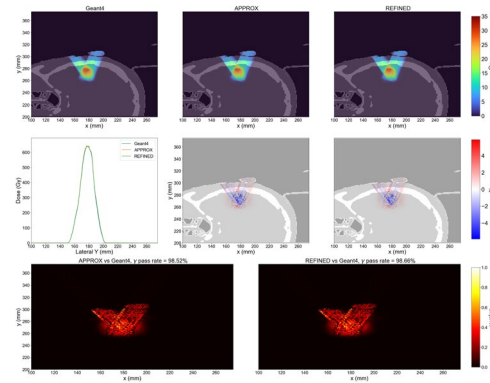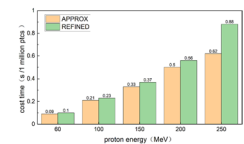


**Fig. 14** (Color online) Runtimes of gMCAP under different energies

programs employing simplified models. It demonstrates strong agreement with Geant4 and achieves a precision similar to that of current GPU-based Monte Carlo codes [29, 30] that utilize simplified models. For homogeneous phantoms, gMCAP achieved a gamma passing rate exceeding 99% for a 2 mm/2% criterion. Moreover, the REFINED model in gMCAP outperformed the APPROX codes by achieving a gamma passing rate exceeding 95% for all phantoms with a 1 mm/1% criterion. These results highlight the superior performance and precision of gMCAP in simulating proton dose distributions.

Additionally, gMCAP demonstrated favorable execution times. For instance, when simulating 10 million primary protons with an energy of 150 MeV in a water phantom consisting of 1 mm voxels, gMCAP required 3.7 s using an NVIDIA RTX 4090. In comparison, the Moqui software [32] required 18 s on an NVIDIA GTX 2080 Ti for the same simulation. Similarly, when simulating 10 million primary protons with an energy of 200 MeV, gMCAP [29] took 5.6 s on an NVIDIA RTX 4090, whereas the gPMC software took 21 s on a Tesla C2050. These results demonstrate that gMCAP offers efficient execution times, outperforms certain alternative software packages, and leverages the power of modern GPUs for accelerated Monte Carlo simulations.

Compared to the APPROX model, the REFINED model in gMCAP improves the gamma pass rate by 3%. This
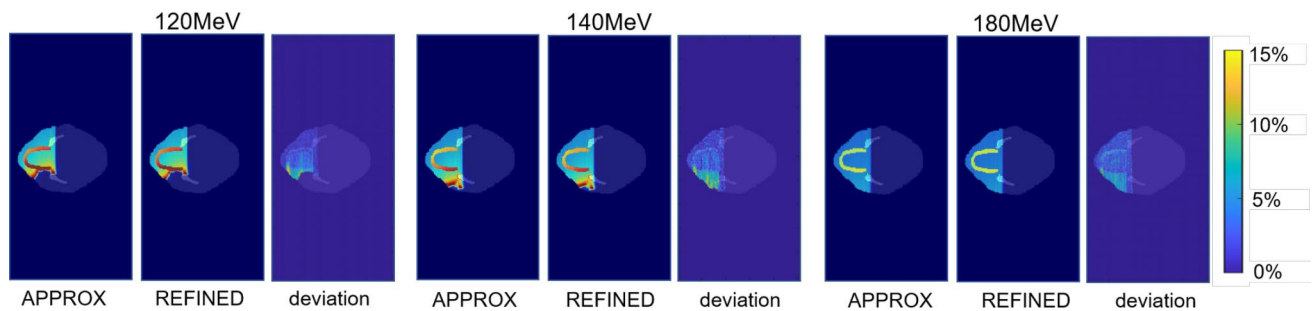
**Fig. 15** (Color online) Comparative analysis of the REFINED and APPROX models employing 120 MeV, 140 MeV, and 180 MeV protons for irradiating the mandibular area, respectively

enhancement is particularly significant in tissue regions rich in elements such as Ca and P, which are abundant in the bones and teeth. In these regions, the difference between the REFINED and APPROX models can reach up to 15%. To ensure calculation accuracy, the REFINED model should be employed in tissue regions containing a higher proportion of middle-$Z$ and high-$Z$ components.

Despite offering a refined nuclear model for precise discrete proton reactions, the current implementation of gMCAP has certain limitations. One is that neutral particles are not considered, which can lead to specific discrepancies in dose calculations. Additionally, the absence of secondary electron transport may affect the evaluation of the radiation dose in the lungs [29].

To address these limitations, we plan to integrate our previous research on electron transport [48] into gMCAP in the next phase of development. Furthermore, we are exploring the integration of gMCAP with the open-source treatment planning system (TPS) matRad for dose calculation. In the future, we envision enhancing gMCAP as a versatile TPS calculation engine. Further improvements and advancements are required to make this system a comprehensive and reliable tool for treatment planning. Please note that the gMCAP code is still in the development stage and is not yet open source.

## Declarations

## References

1. J.L. Chen, S.J. Yun, T.K. Dong et al., Studies of the radiation environment on the Mars surface using the Geant4 toolkit. Nucl. Sci. Tech. **33**, 11 (2022). https://doi.org/10.1007/s41365-022-00987-2

2. S.C. Huang, H. Zhang, K. Bai et al., Monte Carlo study of the neutron ambient dose equivalent at the heavy ion medical machine in Wuwei. Nucl. Sci. Tech. **33**, 119 (2022). https://doi.org/10.1007/s41365-022-01093-z

3. X.Y. Luo, R. Qiu, Z. Wu et al., THUDose PD: a three-dimensional Monte Carlo platform for phantom dose assessment. Nucl. Sci. Tech. **34**, 164 (2023). https://doi.org/10.1007/s41365-023-01315-y

4. H.F. Ou, B. Zhang, S.H. Zhao, Monte Carlo simulation for calculation of fragments produced by 400 MeV/u carbon ion beam in water. Nucl. Instrum. Meth. B **396**, 18–25 (2017). https://doi.org/10.1016/j.nimb.2017.01.077

5. W.Y. Wang, Y.Y. Ma, H. Zhang et al., Comparison between 4D robust optimization methods for carbon-ion treatment planning. Nucl. Sci. Tech. **34**, 139 (2023). https://doi.org/10.1007/s41365-023-01285-1

6. H.F. Ou, B. Zhang, S.J. Zhao, Gate/Geant4-based Monte Carlo simulation for calculation of dose distribution of 400 MeV/u carbon ion beam and fragments in water. Nucl. Sci. Tech. **27**, 83 (2016). https://doi.org/10.1007/s41365-016-0097-3

7. Y. Shi, M.Z. Zhang, L.H. Ou-Yang et al., Design of a rapid-cycling synchrotron for flash proton therapy. Nucl. Sci. Tech. **34**, 145 (2023). https://doi.org/10.1007/s41365-023-01283-3

8. M.Z. Zhang, D.M. Li, L.R. Shen et al., SAPT: a synchrotron-based proton therapy facility in Shanghai. Nucl. Sci. Tech. **34**, 148 (2023). https://doi.org/10.1007/s41365-023-01293-1

9. M.F. Han, J.X. Zheng, X.H. Zeng et al., Investigation of combined degrader for proton facility based on BDSIM/FLUKA Monte Carlo methods. Nucl. Sci. Tech. **33**, 17 (2022). https://doi.org/10.1007/s41365-022-01002-4

10. M. Liu, C.X. Yin, K.C. Chu et al., A scheme design of collimator for gantry in proton therapy facility. Nucl. Instrum. Meth. A **791**, 47–53 (2015). https://doi.org/10.1016/j.nima.2015.04.045

11. Z.Y. Yao, Y.S. Xiao, J.Z. Zhao, Dose reconstruction with Compton camera during proton therapy via subset-driven origin ensemble and double evolutionary algorithm. Nucl. Sci. Tech. **34**, 59 (2023). https://doi.org/10.1007/s41365-023-01207-1

12. L. Deng, G. Li, T. Ye et al., MCDB Monte Carlo Code with Fast Track Technique and Mesh Tally Matrix for BNCT. J. Nucl. Sci. Technol. **44**, 1518–1525 (2007). https://doi.org/10.1080/18811248.2007.9711401

13. C. Wu, Y.H. Pu, X. Zhang, GPU-accelerated scanning path optimization in particle cancer therapy. Nucl. Sci. Tech. **30**, 56 (2019). https://doi.org/10.1007/s41365-019-0582-6

14. H. Paganetti, Range uncertainties in proton therapy and the role of Monte Carlo simulations. Phys. Med. Biol. **57**, R99 (2012). https://doi.org/10.1088/0031-9155/57/11/R99

15. P.L. Petti, Differential pencil beam dose calculations for charged particles. Med. Phys. **19**, 137–149 (1992). https://doi.org/10.1118/1.596887

16. L. Hong, M. Goitein, M. Bucciolini et al., A pencil beam algorithm for proton dose calculations. Phys. Med. Biol. **41**, 1305 (1996). https://doi.org/10.1088/0031-9155/41/8/005

17. M. Soukup, M. Fippel, M. Alber, A pencil beam algorithm for intensity modulated proton therapy derived from Monte Carlo simulations. Phys. Med. Biol. **50**, 5089 (2005). https://doi.org/10.1088/0031-9155/50/21/010

18. J.S. Li, B. Shahine, E. Fourkal et al., A particle track-repeating algorithm for proton beam dose calculation. Phys. Med. Biol. **50**, 1001 (2005). https://doi.org/10.1088/0031-9155/50/5/022

19. P.P. Yepes, D. Mirkovic, P.J. Taddei, A GPU implementation of a track-repeating algorithm for proton radiotherapy dose calculations. Phys. Med. Biol. **55**, 7107 (2010). https://doi.org/10.1088/0031-9155/55/23/S11

20. K.A. Gifford, J.L. Horton, T.A. Wareing et al., Comparison of a finite-element multigroup discrete-ordinates code with Monte Carlo for radiotherapy calculations. Phys. Med. Biol. **51**, 2253 (2006). https://doi.org/10.1088/0031-9155/51/9/010

21. A. Fogliata, G. Nicolini, A. Clivio et al., Accuracy of Acuros XB and AAA dose calculation for small fields with reference to RapidArc® stereotactic treatments. Med. Phys. **38**, 6228–6237 (2011). https://doi.org/10.1118/1.3654739

22. M. Fippel, M. Soukup, A Monte Carlo dose calculation algorithm for proton therapy. Med. Phys. **31**, 2263–2273 (2004). https://doi.org/10.1118/1.1769631

23. F. Salvat, A generic algorithm for Monte Carlo simulation of proton transport. Nucl. Instrum. Meth. B **316**, 144–159 (2013). https://doi.org/10.1016/j.nimb.2013.08.035

24. X. Jia, X. Gu, J. Sempau et al., Development of a GPU-based Monte Carlo dose calculation code for coupled electron-photon transport. Phys. Med. Biol. **55**, 3077 (2010). https://doi.org/10.1088/0031-9155/55/11/006

25. X. Jia, X. Gu, Y.J. Graves et al., GPU-based fast Monte Carlo simulation for radiotherapy dose calculation. Phys. Med. Biol. **56**, 7017 (2011). https://doi.org/10.1088/0031-9155/56/22/002

26. L. Su, Y. Yang, B. Bednarz et al., ARCHERRT-A GPU-based and photon-electron coupled Monte Carlo dose computing engine for radiation therapy: Software development and application to helical tomotherapy. Med. Phys. **41**, 071709 (2014). https://doi.org/10.1118/1.4884229

27. S. Hissoiny, B. Ozell, H. Bouchard et al., GPUMCD: a new GPU-oriented Monte Carlo dose calculation platform. Med. Phys. **38**, 754–764 (2011). https://doi.org/10.1118/1.3539725

28. A.K. Hu, R. Qiu, H. Liu et al., THUBrachy: fast Monte Carlo dose calculation tool accelerated by heterogeneous hardware for high-dose-rate brachytherapy. Nucl. Sci. Tech. **32**, 32 (2021). https://doi.org/10.1007/s41365-021-00866-2

29. X. Jia, J. Schümann, H. Paganetti et al., GPU-based fast Monte Carlo dose calculation for proton therapy. Phys. Med. Biol. **57**, 7783 (2012). https://doi.org/10.1088/0031-9155/57/23/7783

30. A. Schiavi, M. Senzacqua, S. Pioli et al., Fred: a GPU-accelerated fast-Monte Carlo code for rapid treatment plan recalculation in ion beam therapy. Phys. Med. Biol. **62**, 7482 (2017). https://doi.org/10.1088/1361-6560/aa8134

31. D. Maneval, B. Ozell, P. Després, pGPUMCD: an efficient GPU-based Monte Carlo code for accurate proton dose calculations. Phys. Med. Biol. **64**, 085018 (2019). https://doi.org/10.1088/1361-6560/ab0db5

32. H. Lee, J. Shin, J.M. Verburg et al., MOQUI: an open-source GPU-based Monte Carlo code for proton dose calculation with efficient data structure. Phys. Med. Biol. **67**, 174001 (2022). https://doi.org/10.1088/1361-6560/ac8716

33. W.G. Li, C. Chang, Y. Qin et al., GPU-based cross-platform Monte Carlo proton dose calculation engine in the framework of Taichi. Nucl. Sci. Tech. **34**, 77 (2023). https://doi.org/10.1007/s41365-023-01218-y

34. D.R. White, R.V. Griffith, I.J. Wilson, The composition of body tissues. Rep. Int. Commiss. Radiat. Units Meas. **os–24**, 5–9 (1992). https://doi.org/10.1093/jicru_os24.1.5

35. H.H. Barschall, M.B. Chadwick, D.T.L. Jones et al., Appendix D: proton data tables. Rep. Int. Commiss. Radiat. Units Meas. **os–32**, 166–236 (2000). https://doi.org/10.1093/jicru_os32.2.166

36. N. Qin, P. Botas, D. Giantsoudi et al., Recent developments and comprehensive evaluations of a GPU-based Monte Carlo package for proton therapy. Phys. Med. Biol. **61**, 7347 (2016). https://doi.org/10.1088/0031-9155/61/20/7347

37. Basic anatomical and physiological data, The skeleton. Ann. ICRP **25**, 1–80 (1995). https://doi.org/10.1016/S0146-6453(00)80004-4

38. M.B. Chadwick, P.M. DeLuca, R.C. Haight, Nuclear data needs for neutron therapy and radiation protection. Radiat. Port. Dosim **70**, 1–12 (1997). https://doi.org/10.1093/oxfordjournals.rpd.a031920

39. Y. Fu, Y. Lei, T. Wang et al., A review of deep learning based methods for medical image multi-organ segmentation. Phys. Medica **85**, 107–122 (2021)

40. U. Schneider, E. Pedroni, A. Lomax, The calibration of CT Hounsfield units for radiotherapy treatment planning. Phys. Med. Biol. **41**, 111 (1996). https://doi.org/10.1088/0031-9155/41/1/009

41. H.H.C. Lee, Y.K. Park, X. Duan et al., Convolutional neural network based proton stopping-power-ratio estimation with dual-energy CT: a feasibility study. Phys. Med. Biol. **65**, 215016 (2020). https://doi.org/10.1088/1361-6560/abab57

42. R. Zhang, W.D. Newhauser, Calculation of water equivalent thickness of materials of arbitrary density, elemental composition and thickness in proton beam irradiation. Phys. Med. Biol. **54**, 1383 (2009). https://doi.org/10.1088/0031-9155/54/6/001

43. S. Agostinelli, J. Allison, K. Amako et al., GEANT4-a simulation toolkit. Nucl. Instrum. Meth. A **506**, 250–303 (2003). https://doi.org/10.1016/S0168-9002(03)01368-8

44. B. Rossi, K. Greisen, Cosmic-ray theory. Rev. Mod. Phys. **13**, 240 (1941). https://doi.org/10.1103/RevModPhys.13.240

45. J. Ranft, Estimation of radiation problems around high-energy accelerators using calculations of the hadronic cascade in matter. Part. Accel. **3**, 129–161 (1972). https://doi.org/10.1130/0016-7606(1981)92$<$197:TPODRA$>$2.0.CO;2

46. H.P. Wieser, E. Cisternas, N. Wahl et al., Development of the open-source dose calculation and optimization toolkit matRad. Med. Phys. **44**, 2556–2568 (2017). https://doi.org/10.1002/mp.12251

47. D.A. Low, W.B. Harms, S. Mutic et al., A technique for the quantitative evaluation of dose distributions. Med. Phys. **25**, 656–661 (1998). https://doi.org/10.1118/1.598248

48. Z. Wu, W. Lu, S.C. Yan et al., Development of a GPU-accelerated photon-electron coupled transportation Monte Carlo code and its application. J. Harb. Eng. Uni **43**, 1649–1656 (2022). https://doi.org/10.11990/jheu.202206028