

# Research on Brain Tumor Segmentation Deployment Method Based on Mixed Precision Inference Acceleration

Yamin Wang, Beibei Hou

School of Computer Science and Technology, Henan Polytechnic University, Jiaozuo 454000, China

---

**Abstract:** Aiming at the engineering bottlenecks faced by deep learning brain tumor segmentation models during real-world clinical deployment, such as high inference latency, large peak GPU memory consumption, and redundant model file sizes, this paper proposes a Mixed Precision Acceleration (MPA) inference deployment method. This method exports the standard FP32 precision model via the ONNX format and relies on the NVIDIA TensorRT inference engine to perform deep optimizations in FP16 half-precision mode, including computational graph redundancy elimination, layer fusion, kernel auto-tuning, and dynamic memory reuse. Experiments under three representative missing-modality scenarios on the BraTS2020 dataset demonstrate that, under the premise of maintaining nearly lossless segmentation accuracy (evaluation metric differences are all less than 0.01%), the MPA method compresses the model file volume by 70.6% on average (reduced to 17.5MB), significantly decreases peak memory consumption by 82.6%~85.7% (requiring only 700MB at minimum), and improves inference speed by 13.6%~16.5%. Compared to the Torch-Compile compilation scheme, the MPA method achieves a more optimal trade-off among inference speed, memory overhead, and storage cost, providing a highly feasible engineering solution for deploying brain tumor segmentation models to computationally constrained environments such as primary medical workstations, mobile terminals, and intraoperative navigation systems.

**Keywords:** Brain tumor segmentation; Mixed Precision Acceleration; TensorRT.

---

## 1. Introduction

In recent years, deep learning technology has made breakthrough progress in multi-modality magnetic resonance brain tumor image segmentation tasks [1-5], significantly improving the accuracy of automatic delineation of target lesion regions [6-10]. However, most of these high-precision deep learning networks [11-13] are developed and validated in laboratory environments with abundant computing resources [14,15]. Accompanying the improvement in algorithm precision, the computational complexity and parameter volume of the models also undergo a dramatic increase [16,17]. A significant engineering gap still exists between model prototypes in laboratory settings and resource-constrained clinical terminal deployments [18-20]. An excellent deep learning segmentation model requires not only theoretical accuracy guarantees but also the ability to overcome performance constraints in actual clinical applications [21,22].

When pushing a fully-trained network into a real clinical deployment environment, three major engineering bottlenecks arise [23]: First, the standard FP32 single-precision floating-point inference under the PyTorch framework requires constructing a complete dynamic computational graph at runtime and unconditionally allocating memory operation by operation. The inference latency for a single sample ranges from 3.27 to 3.36 seconds, making it difficult to meet the needs of clinical scenarios requiring near real-time feedback, such as intraoperative navigation [24]. Second, the peak GPU memory consumption during inference reaches as high as 4.88~5.02GB, imposing strict hardware constraints on primary medical workstations or mobile medical terminals equipped with only entry-level GPUs. Third, the model file stored in FP32 format is about

59.5MB in size, which increases the time cost of model transmission and loading in bandwidth-limited remote diagnosis or edge computing scenarios [25].

To address the core pain points at the deployment level mentioned above, this section proposes a Mixed Precision Acceleration (MPA) inference optimization method. The technical route of this method is to export the trained FP32 precision network into the Open Neural Network Exchange (ONNX) intermediate representation format and subsequently utilize the NVIDIA TensorRT high-performance inference engine to execute deep optimization on the computational graph in FP16 half-precision mode. During the process of building the inference engine, TensorRT sequentially performs a series of optimization operations: redundant computation node elimination, layer fusion of adjacent operators, optimal CUDA kernel auto-tuning oriented to specific GPU architectures, and dynamic tensor memory reuse during inference. By degrading most computation-intensive convolution and matrix operations in the network from FP32 to FP16, the MPA method fully unleashes the native hardware acceleration capability of the Tensor Cores in modern NVIDIA GPUs for half-precision floating-point operations. On the Ampere architecture of the NVIDIA RTX 3090 GPU, its third-generation Tensor Cores provide twice the theoretical FP16 computing power as FP32, and FP16 data occupies only half the memory bandwidth of FP32, thereby obtaining acceleration benefits in both computational throughput and memory access dimensions.

The main content of this study is arranged as follows: Section 2 systematically illustrates the key technical principles upon which mixed precision inference acceleration relies, including the floating-point precision system, half-precision quantization mechanism, core optimization technologies of the TensorRT inference engine, and the end-

to-end mixed precision acceleration inference pipeline designed in this paper. Section 3 comprehensively compares the MPA method with two baseline schemes—native PyTorch inference and Torch-Compile compilation optimization—under three representative missing-modality scenarios on the BraTS2020 dataset, systematically verifying the effectiveness of the MPA method from four dimensions: segmentation accuracy retention, inference speed improvement, memory consumption reduction, and model storage compression. Section 4 summarizes the work of this paper.

## 2. Mixed Precision Inference Acceleration Technology

This section systematically elaborates on the key technical principles on which the mixed precision inference acceleration method relies. First, starting from the underlying mechanism of floating-point representation, it analyzes the information coding differences between FP32 and FP16 and their impacts on model inference. Second, it deeply introduces the core optimization technology stack of the NVIDIA TensorRT inference engine. Finally, based on the above theoretical foundation, it presents the end-to-end mixed precision accelerated inference workflow designed in this paper.

### 2.1. Floating-point Precision and Half-precision Quantization Principles

In the representation and computation of deep learning models, floating-point precision directly determines the numeric dynamic range, calculation accuracy, and storage bandwidth requirements. The two most commonly used floating-point formats in current deep learning training and inference are the FP32 single-precision floating-point and FP16 half-precision floating-point defined by the IEEE 754 standard, which inherently differ in bit-width encoding.

FP32 utilizes 32-bit binary encoding, consisting of 1 sign bit, 8 exponent bits, and 23 mantissa bits. It can represent a dynamic range from approximately  $1.18 \times 10^{-38}$  to  $3.40 \times 10^{38}$ , with an effective decimal precision of about 7 digits. Its value is determined by the following formula:

$$\text{Value}_{\text{FP32}} = (-1)^s \times 2^{(e-127)} \times \left( 1 + \sum_{i=1}^{23} b_{23-i} \cdot 2^{-i} \right)$$

where  $s$  is the sign bit,  $e$  is the unsigned integer representation of the exponent with a bias of 127, and  $b_i$  is the binary value of the  $i$ -th bit of the mantissa.

FP16 employs 16-bit binary encoding, consisting of 1 sign bit, 5 exponent bits, and 10 mantissa bits. It has a dynamic range of about  $6.10 \times 10^{-5}$  to  $6.55 \times 10^4$  and an effective precision of 3 to 4 decimal digits. Its calculation formula is:

$$\text{Value}_{\text{FP16}} = (-1)^s \times 2^{(e-15)} \times \left( 1 + \sum_{i=1}^{10} b_{10-i} \cdot 2^{-i} \right)$$

where the exponent bias is 15. The machine precision for the two formats is  $\epsilon_{\text{FP32}} = 2^{-23} \approx 1.19 \times 10^{-7}$  and  $\epsilon_{\text{FP16}} = 2^{-10} \approx 9.77 \times 10^{-4}$ , respectively, with a ratio of upper bounds for a single rounding error reaching approximately 4096 times.

The essence of the quantization from FP32 to FP16 is lossy compression: the mantissa bits shrink from 23 to 10 bits, causing the decimal parts precision to drop by about three orders of magnitude; the exponent bits shrink from 8 to 5 bits, leading to a substantial narrowing of the dynamic range.

However, in neural network inference, the distribution of weights and intermediate activations usually concentrates within a relatively small numerical range, and the final voxel-level segmentation decision is a statistical aggregation result of massive floating-point multiply-add operations. The approximate  $O(10^{-4})$  magnitude error introduced by a single FP16 operation cancels out due to sign randomness over millions of accumulations, leading to an impact on the overall prediction result that approaches zero. This provides a solid theoretical foundation for adopting a mixed-precision strategy in half-precision inference.

From an execution hardware perspective, the efficiency improvements of FP16 inference manifest at three levels:

Firstly, the doubling effect of storage bandwidth. The bit-width of FP16 is precisely half that of FP32. Accordingly, the number of bytes transferred between the GPUs global memory and the SM cache halves. For bandwidth-intensive operators like 3D convolution, theoretically, the number of elements transferred per memory read and write doubles, fundamentally alleviating the memory bandwidth bottleneck.

Secondly, native acceleration by Tensor Cores. Starting from the Volta architecture, NVIDIA introduced Tensor Cores specifically designed for mixed-precision matrix multiplication in GPUs. Taking the RTX 3090 (Ampere architecture, third-generation Tensor Core) used in experiments as an example, its FP16 matrix multiply-add throughput reaches 156 TFLOPS, far exceeding the 35.6 TFLOPS of FP32 CUDA Cores. When generating an FP16 inference plan, TensorRT maps convolution operations meeting size alignment requirements directly to the HMMA instruction path of Tensor Cores. The Ampere architecture uses a  $16 \times 16 \times 16$  matrix block as the basic computational unit, completing one matrix block multiply-add per clock cycle:

$$D = A \times B + C$$

Here,  $A$  and  $B$  are loaded in FP16 format, while the accumulated results  $C$  and  $D$  remain at FP32 precision. This design not only doubles the computational power via FP16 but also prevents precision overflow during large-scale summation via FP32 intermediate accumulation.

Thirdly, significant reduction in GPU memory footprint. After weights are stored in FP16, the model file volume halves directly. Furthermore, intermediate activation tensors of each layer also reside in memory in FP16. In cooperation with TensorRT's static memory reuse planning, the peak memory reduction is usually far greater than the theoretical 50% lower bound. According to the measurements in this study's experiments, this achieves an 82% to 86% reduction, which is vital for deploying models on entry-level GPU-equipped medical workstations or embedded medical terminals.

### 2.2. TensorRT Inference Engine Optimization Technologies

NVIDIA TensorRT is a high-performance deep learning inference engine aimed at deployment in production environments [27]. Its ultimate core goal is to maximize the utilization of GPU hardware parallelism and compress the inference latency to an engineering limit. Its optimization workflow is split into a build phase and an execution phase. The key technologies involved in this paper consist of the following four main aspects.

First, computational graph optimization and redundancy elimination. After reading the ONNX format network description [26], TensorRT initially applies static graph

transformations on the native graph: removing redundant nodes such as identity mappings and invalid Reshapes; executing constant folding on subgraphs that can be determined at compile time to eliminate runtime repeated calculations; and identifying cascaded Conv-GN-ReLU structures to merge them into a single computing kernel, which offsets multiple memory reads/writes for intermediate feature maps. These transformations substantially lower runtime operator scheduling times and memory access frequencies without altering the forward propagation semantics.

Second, layer fusion and operator fusion. Layer fusion is TensorRT's core mechanism for realizing high-efficiency inference. It combines logically independent but data-flow-tightly-coupled multi-operators into a single composite operator, finishing multi-step calculations within a single CUDA kernel call. Taking the 3D Convolution  $\rightarrow$  Group Normalization  $\rightarrow$  ReLU triplet widely used in this network as an example: before fusion, each of the three operators requires one global memory read and one global memory write, totaling six memory accesses. However, after being fused into a unified 'ConvGNReLU' kernel, the input tensor is only queried once from global memory, and the final result is written back once, decreasing memory access from 6 to 2. This drops global memory bandwidth consumption by 66.7%. Extending this to a network featuring  $L$  fusible triplets, the ratio of total memory accesses before and after fusion represents:

$$\frac{M_{\text{fused}}}{M_{\text{unfused}}} = \frac{2L}{6L} = \frac{1}{3}$$

Beyond saving memory accesses, layer fusion prevents fixed startup overhead for multiple kernel scheduling. Since each CUDA kernel scheduling takes tens of microseconds, the cumulative benefit is highly observable for deep networks.

Third, kernel auto-tuning. During the build phase, TensorRT performs real-time profiling on the target GPU for each network layer. It screens out the lowest-latency implementation scheme from multiple built-in choices. For identical convolution operations, candidate methods involve cuBLAS GEMM implementation via 'im2col' matrix unfolding, manually tuned cuDNN convolution kernels, or HMMA-level kernel functions tailored for Tensor Cores [28]. After auto-tuning, the engine permanently maps the optimal strategy into serialized binary executable files. No re-tuning is required at runtime, making certain it executes following optimal paths under varying input batches.

Fourth, dynamic tensor memory management. In the standard PyTorch dynamic graph mode, each operator independently allocates and releases memory. Frequent memory allocations not only expand system call overhead but additionally induce memory fragmentation. TensorRT replaces this with static memory planning: in the build phase, it examines tensor life cycles across the entire network, mapping tensor pairs non-conflicting in the temporal dimension into overlapping physical memory blocks to enact reuse. Assuming a network retains  $N$  intermediate activation tensors featuring volumes  $S_1, S_2, \dots, S_N$ , TensorRT sorts them into multiple overlap-free reuse groups applying interval coloring algorithms. The requisite memory volume totals merely the sum of the maximum tensor dimensions embedded in each group, differing vastly from the aggregation of all activations. Furthermore, all intermediate computations wrap up within a pre-allocated fixed-size workspace, unconditionally erasing the runtime latency caused by

dynamic memory governance. This represents the underlying reason our MPA method achieves peak memory reductions outperforming the 50% theoretical ceiling presented by FP16 quantization [30].

## 2.3. Mixed Precision Acceleration Inference Workflow Design

Based on the technologies illustrated above, this paper designs an end-to-end mixed precision accelerated inference pipeline [29]. It transforms the initially trained FP32 precision network into deeply optimized FP16 TensorRT inference engines. The comprehensive process scales across four stages.

### Phase 1: Model Export.

Utilize PyTorch's ONNX export API to serialize the deep network into a framework-agnostic ONNX arrangement. The file completely encodes the overarching forward calculation logic driving from multi-channel MRI inputs to the three-region segmentation output, including all 3D convolution layers, group normalization layers, activation functions, trilinear upsampling, as well as tensor concatenation transactions found at skip connections. The input tensor dimensions are permanently anchored during export as  $(N, C, D, H, W)$ , wherein  $N$  marks the batch size,  $C$  depicts input channels (fetching 1/2/3 responding to multiple modality missing contexts), and  $D \times H \times W$  indicates voxel resolution. Moreover, a specific opset edition is established to certify parsing compatibility throughout TensorRT.

### Phase 2: TensorRT Engine Construction.

After reading the ONNX file, instantiate a 'BuilderConfig' assigning the FP16 precision indicator to initialize deep optimization. It should be declared that TensorRT functioning under FP16 mode rarely delegates all layers to half-precision blindly. Through built-in accuracy calibration logic, it deciphers the exceptional scheme layer-by-layer: numeric-sensitive normalization processes alongside element-wise additions in residual connections remain fixed at FP32 to inhibit overflow. Contrastingly, computation-dense wide-dimensional convolution operations receive prioritized direction toward the FP16 Tensor Core routes. Such an adaptive layer-by-layer selection fundamentally epitomizes the core definition of "Mixed Precision." Following construction, the engine maps to binary sequences encapsulating FP16-configured weights accompanied by explicitly compiled CUDA kernels.

### Phase 3: Inference Execution.

Post-deserializing the constructed engine representation, trigger inference invoking the CUDA streams asynchronous interfaces. The triplet actions covering data bridging (H2D), forward estimations managed by the GPU, and resulting retrievals (D2H), are asynchronously scheduled following strict queuing within identical CUDA streams. Discarding python-layer operator-by-operator delays, this strategy exclusively answers to the CUDA runtime environment.

### Phase 4: Accuracy Validation and Consistency Auditing.

To ascertain that the quantization and acceleration phase inserts zero clinically unacceptable segmentation discrepancies, experimental checkpoints enforce value comparisons. Having generated the TensorRT engine, the evaluation cross-references outcomes produced by PyTorch FP32 logic juxtaposed with the TensorRT FP16 algorithm systematically iterating over the complete test set composed of 369 validation specimens. Specifically indexing the Whole Tumor (WT), Tumor Core (TC), and Enhancing Tumor (ET) regional zones, multi-faceted index differences spanning Dice

similarity coefficients (DSC), Hausdorff distance (HD95), and Sensitivities authenticate mixing-precision retaining limits.

### 3. Experimental Results and Analysis

#### 3.1. Experimental Setup

The operational hardware platform matching standard software orchestration parallels earlier configurations, utilizing an Ubuntu 20.04 system equipped with an NVIDIA GeForce RTX 3090 containing 24GB of memory performing total assessments. NVIDIA TensorRT variation 8.x, precisely complementing the matching CUDA toolkit layout, supports engine processing. ONNX formulation operates fixedly upon opset volume 17.

In order to sufficiently estimate generalized performance yielded by the presented MPA methodology targeting fluctuating scale reductions reflecting clinical reality, this chapter nominates three distinctive modality composition variations representing varied degradation contexts: (1) Triple-modality combination enclosing FLAIR+T1ce+T2 (characterized by 3-channel input), capturing standard situations absent of independent T1w modality configurations; (2) Dual-modality combination aggregating T1ce+T2 (2-channel input), portraying intermediate degradations lacking dual elements; (3) Uni-modality T2 setup (1-channel input), symbolizing extreme limits restricted to simply foundational sequence mapping. These combinations encircle maximum spectrum bandwidth defining input accessibility ranging from moderate availability down through stark scarcity.

Establishing contrastive observation schemas evaluating inference proficiency involves structuring three distinctive benchmark methodologies extending cross-comparisons:

Baseline 1: Native PyTorch Inference. Proceeding faithfully traversing initial configurations applying straight FP32 precision strictly within default 'eval' operating constraints, stripping supplemental optimizations, serving as primary baselines.

Baseline 2: MPA Mixed Precision Acceleration. Signifying the deployment algorithm endorsed in this chapter configuring pre-computed models exported towards ONNX interpretations routed across TensorRT frameworks executing deeply fine-tuned FP16 logic algorithms; noticeably tightening parameter capacities concurrently decreasing memory encumbrances.

Baseline 3: Torch-Compile Compilation Optimization. Harvesting functional compilation advancements implemented by PyTorch 2.x directing instantaneous graph assessments optimizing operator mergers alongside target hardware CUDA instruction formations. Maintaining unbroken FP32 properties yielding boosted execution velocity profiles signifies supreme PyTorch-ecosystem representations avoiding structural manipulation.

Individually confronting singular approaches handling varied modality organizations invokes iterative sample calculations reviewing the fully composed BraTS2020 testing assortment extending through 369 records. Evaluators document specific metrics reflecting inference duration per file (restricted towards frontward calculation propagations), inclusive processing timeframe (packaging preliminary dataset handling followed by concluding filtering calculations), and overarching absolute GPU consumption (totalizing distributed memory combined adjacent allocated reserves) incorporating fundamental blueprint memory size dimensions. Evaluational records synchronously analyze Dice outcomes matching alongside equivalent HD95 ranges accompanied by descriptive Sensitivities dissecting respective WT, TC, and ET subregional segmentations. All temporal indices reflect statistical averages spanning the 369 sequences, eliminating localized distortions.

#### 3.2. Comparison Experiments

Table 1 details the comparative results of segmentation accuracy for the three inference schemes under varying modality combinations.

**Table 1.** Comparison of segmentation accuracy under different combinations of modalities using various reasoning methods

| Modality combination | Inference methods | Dice (%) |        |        | HD95(mm) |        |        | Sensitivity (%) |        |        |
|----------------------|-------------------|----------|--------|--------|----------|--------|--------|-----------------|--------|--------|
|                      |                   | WT       | TC     | ET     | WT       | TC     | ET     | WT              | TC     | ET     |
| FLAIR<br>+T1ce+T2w   | PyTorch           | 90.774   | 85.613 | 77.938 | 6.095    | 5.902  | 5.127  | 89.589          | 81.691 | 79.425 |
|                      | MPA               | 90.774   | 85.612 | 77.935 | 6.098    | 5.905  | 5.130  | 89.586          | 81.691 | 79.423 |
|                      | Torch-Compile     | 90.774   | 85.613 | 77.938 | 6.053    | 5.902  | 5.127  | 89.589          | 81.691 | 79.425 |
| T1ce+T2w             | PyTorch           | 87.693   | 85.116 | 76.477 | 8.750    | 5.832  | 5.226  | 86.391          | 83.725 | 79.885 |
|                      | MPA               | 87.695   | 85.119 | 76.478 | 8.746    | 5.831  | 5.227  | 86.400          | 83.736 | 79.892 |
|                      | Torch-Compile     | 87.693   | 85.116 | 76.476 | 8.750    | 5.832  | 5.226  | 86.391          | 83.725 | 79.885 |
| T2w                  | PyTorch           | 86.868   | 68.777 | 42.268 | 11.276   | 12.536 | 11.572 | 86.737          | 74.149 | 50.070 |
|                      | MPA               | 86.868   | 68.775 | 42.261 | 11.277   | 12.533 | 11.570 | 86.725          | 74.150 | 50.065 |
|                      | Torch-Compile     | 86.868   | 68.777 | 42.268 | 11.276   | 12.536 | 11.572 | 86.737          | 74.149 | 50.070 |

Deriving crucial considerations extracted from Table 1 indicates central phenomena:

(1) The applied MPA structure continuously manifests strictly non-degrading representation boundaries managing

universal testing inputs capturing diverse subregions. Concerning triple arrangement settings dictating FLAIR+T1ce+T2 availability, index metrics framing PyTorch Native methodologies directly coordinate MPA

conclusions matching DSC parameters aligned seamlessly across secondary rounding calculations (DSC scores identifying WT/TC/ET read strictly maintaining 90.77%/85.61%/77.94%). Reflecting Dual-modality boundaries corresponding up T1ce+T2 provisions, statistical representations yield sporadic negligible upwards bias evaluating MPA-directed TC Sensitivity scoring recording fractional surges bounding 0.02% (83.74% contrasting 83.72%). These observations derive foundational origins targeting FP16 random distribution variations originating stochastically encouraging minimal forward drifts consistently retaining clinical standards. Investigating extreme isolation involving sole-channel operations mapping structural T2 environments, comparison highlights minute discrepancies tracking roughly limited near 0.01% boundaries measuring ET classifications (PyTorch ET indexing 42.26% opposed against baseline expectations scoring 42.27%).

These experimental results fully demonstrate that the FP16 mixed-precision quantization strategy in this paper has achieved lossless inference compression in brain tumor segmentation tasks in a true sense. From the perspective of numerical analysis, the massive floating-point multiply-accumulate operations in 3D convolution have a natural error-averaging effect. The relative error of about  $1.22 \times 10^{-4}$  introduced by a single FP16 multiplication acts randomly and is statistically offset after hundreds of thousands of accumulation steps. Therefore, the impact on the final voxel-level segmentation result approaches zero.

(2) The Torch-Compile compilation architecture also upholds firmly parallel precision levels resembling initial unoptimized inference because the architecture consistently preserves floating-point numerical accuracy securely matching default FP32 conditions, purely emphasizing execution cycle augmentations applying systematic algorithm interpretations. Highly cohesive outcomes intersecting tripartite strategies emphasize our applied speed improvement modifications securely avoiding compromising performance variables highlighting indispensable assurances shielding clinical security margins.

(3) Based on cross-modality precision trends, all three inference architectures display absolutely concurrent behaviors: aligned together with proliferating accessible sequences covering multimodal origins, performance

evaluations matching WT, TC, and ET subregions uniformly enhance relative Dice evaluation marks matching synchronous HD95 downturns. Specifying distinct segments records expanding scores across singular down toward comprehensive models reading WT regions identifying 86.87%, 87.69%, up pointing to 90.77% clarifying FLAIR configuration significance drawing expansive bounds determining surrounding edemas. Correspondingly, tracking ET subregion evaluations traces leaps spanning basic configurations logging 42.26% rocketing vertically confronting triple-element limits capturing 77.94% showcasing undisputed foundational dependency isolating critical contrasting values emphasizing indispensable diagnostic criteria attributed uniquely spanning T1ce sequences highlighting structural representations. These patterns perfectly align with theoretical analyses of the contribution levels of each modality to different tumor subregions, further verifying that the MPA method strictly maintains the multi-modal feature utilization capability of the original model. It did not destroy the networks information fusion logic due to precision quantization. Its especially worth pointing out that under extreme modality deficiency involving mere single T2 elements, performance metrics matching ET representations plunge measuring uniquely 42.26%, because the strong signaling responses characterizing targeted enhancing architectures massively correspond toward comparative shading contrasts delivered universally representing T1ce frameworks highlighting contrast enhancement properties which simply evade capture implementing standalone sequential parameters, rendering elemental network breakdowns limiting precision. This observation is fundamentally constrained by network capabilities and further confirms current mapping relationships concerning underlying multimodal missing information paradigms.

### 3.3. Inference Efficiency and Deployment Analysis

Table 2 comprehensively compares the deployment efficiency metrics of the three inference methods in three dimensions: model file volume, inference speed, and GPU memory consumption.

**Table 2.** Comparison of reasoning efficiency among different reasoning methods

| Modality combination | Inference methods | Model size (MB) | Average inference time (s) | Inference speedup ratio | Peak allocated memory (MB) | Memory reduction rate | Peak reserved memory (MB) |
|----------------------|-------------------|-----------------|----------------------------|-------------------------|----------------------------|-----------------------|---------------------------|
| FLAIR<br>+T1ce+T2w   | PyTorch           | 59.56           | 3.361                      | 1.000×                  | 5018                       | —                     | 5922                      |
|                      | MPA               | 17.54           | 2.885                      | 1.165×                  | 872                        | 82.6%↓                | 997                       |
|                      | Torch-Compile     | 59.56           | 2.571                      | 1.307×                  | 3453                       | 31.2%↓                | 4372                      |
| T1ce+T2w             | PyTorch           | 59.54           | 3.274                      | 1.000×                  | 4949                       | —                     | 5872                      |
|                      | MPA               | 17.46           | 2.882                      | 1.136×                  | 786                        | 84.1%↓                | 863                       |
|                      | Torch-Compile     | 59.54           | 2.482                      | 1.319×                  | 3331                       | 32.7%↓                | 4646                      |
| T2w                  | PyTorch           | 59.53           | 3.349                      | 1.000×                  | 4882                       | —                     | 6250                      |
|                      | MPA               | 17.53           | 2.894                      | 1.157×                  | 700                        | 85.7%↓                | 740                       |
|                      | Torch-Compile     | 59.53           | 2.555                      | 1.311×                  | 3264                       | 33.1%↓                | 4124                      |

Based on the experimental data in Table 2, an in-depth analysis can be conducted from the following dimensions:

(1) Model File Volume Compression Analysis

The MPA method impressively compresses model volume size under the three modality associations moving downwards from baseline 59.56MB, 59.54MB, and 59.53MB scaling

directly arriving alongside newly compressed sizes featuring 17.54MB, 17.46MB, corresponding next to 17.53MB respectively, representing 70.6% gross volume depreciation metrics achieving extraordinary storage capabilities. These noticeable compression results derive mainly scaling spanning dual independent origins: Primarily targeting initial matrix elements transitioning natively exchanging standard FP32 formats moving towards structured reduced FP16 parameters natively ensuring baseline 50% diminutions; Moreover, serialization sequencing structured across TensorRT protocols structurally realigns topological structures streamlining array orientations systematically nullifying underlying redundant structural diagrammatic indicators eradicating leftover systemic variables compounding structural sizes further down.

Notably characterizing negligible variation shifting along diverse independent modality streams (ranging simply covering scales 17.46~17.54MB) reinforces perspectives underlying initial theoretical estimates confirming spatial weights fundamentally track extensive network convolutions circumventing fractional interface parameters associated directly processing localized inputs. Conversing these achievements confronting Torch-Compile arrangements illuminates non-existent file compression returns preserving intact configurations identical pointing towards FP32 baselines retaining full 59.5MB architectures. Essentially interpreting that optimizations native identifying instantaneous Torch-Compile parameters operate explicitly running compilation cycles ignoring hard storage restructuring modifications. For remote telemetry environments alongside mobile peripheral configurations, logging 70.6% magnitude depreciations profoundly economizes communication bandwidth facilitating network configurations updating logic deployments.

### (2) Inference Speed Analysis

Investigating computing speeds traces acceleration performance executing MPA calculations delivering continuous improvements marking 1.165 $\times$ , 1.136 $\times$  along 1.157 $\times$  improvements charting concurrent modality settings driving standard aggregated increases scaling 1.153 $\times$  margins representing roughly an incremental gain logging upwards 15.3% accelerated processing bounds. Inference enhancements derive foundationally combining multiple systematic factors merging smoothly: initially reflecting straight capabilities yielded channeling FP16 computations traversing designated Tensor Cores pushing structural compute levels higher; merging contiguous processing actions implementing Layer Fusion mechanisms functionally negating latency delays spanning iterative GPU kernel invocations eliminating repetitive intermediary reads; ultimately crowning automated profile adjustments selectively mapping most advanced architectural calculation routes executing independent layers tracking maximum throughput parameters.

Comparatively tracking identical tasks executed running independent Torch-Compile operations produces marginally accelerated calculations reading slightly higher speeds registering 1.307 $\times$ , 1.319 $\times$  coupled finally against 1.311 $\times$  registering larger speed jumps defining aggregated averages representing 1.312 $\times$  metric boosts translating matching limits crossing roughly around 31.2% calculation advantages significantly bypassing concurrent MPA computations explicitly tracing unadulterated computing evaluations. Surpassing velocities generated configuring Torch-Compile

procedures associate fundamentally targeting broad scale diagram transformations rendering broad integration linking specific custom native machine CUDA sequences concurrently tapping ample resource structures built intrinsically tracking expansive standalone FP32 CUDA cores mapping natively corresponding heavy RTX 3090 frameworks enabling advanced efficiency calculations mapping deep extensive geometric 3D dimensional processing patterns seamlessly executing expanded calculation ranges. However, as demonstrated below, Torch-Compiles speed advantage is exchanged for a massive memory footprint, limiting its applicability in memory-constrained deployment scenarios.

### (3) GPU Memory Footprint Analysis

Memory occupancy represents the most critical determining variable addressing viability evaluations mapping deployment strategies constrained tracking minimal hardware resources. Documented statistical measurements starkly reinforce overwhelming efficiencies executed strictly traversing MPA techniques managing constrained storage. Figure 1 shows the peak allocated VRAM usage.

In the tri-modal FLAIR+T1ce+T2 scene, the peak allocated memory for PyTorch native inference reaches as high as 5018MB, and reserved memory peaks at 5922MB, nearly occupying a quarter of the RTX 3090s 24GB VRAM. The MPA method drastically scales down the peak allocated memory to 872MB (a reduction of 82.6%), and reserved memory to 997MB (a reduction of 83.2%). This improvement owes to TensorRTs static memory planning and tensor reuse mechanism. By analyzing the life cycle of each memory tensor belonging to the network, TensorRT maps intermediately processed memory arrays featuring zero temporal conflicts pointing uniquely merging inside unified hardware coordinates fully abolishing the inherent memory fragmentation complexities natively affecting standardized PyTorch dynamic evaluations.

Moving along the simulated dual combination T1ce+T2 structures evaluates parallel calculations reflecting MPA allocating simply 786MB boundaries reflecting further improvements representing roughly an incredible 84.1% memory saving dropping reserved structures simultaneously yielding similarly down around baseline limits tracing roughly logging measurements defining just 863MB scales. Modalities lacking singular sequences concurrently trim channel allocations dictating upfront convolutions dropping subsequent volume dimensions mirroring lighter footprints.

Within extreme operational isolation modeling exclusive single parameters spanning independent T2 modalities, evaluated peak consumed dimensions utilizing applied MPA constraints track extraordinarily constrained profiles mapping specifically capturing measurements recording limited minimal peaks mapping unthinkable bare maximums peaking solely spanning measurements scaling logging completely restricted allocations reflecting simply 700MB values recording unprecedented 85.7% bounds alongside comparable measurements mapping barely reserved parameters capturing similarly bounded magnitudes restricting memory down targeting limited minimal 740MB scopes matching extraordinary 88.2% drop representations explicitly compared running baseline equivalent testing structures indicating almost seven mapping towards an incredible full eight magnitudes multiplied improvement.

While mapping operational benefits generated highlighting Torch architecture reflect greater calculation improvements,

tracking subsequent concurrent memory expenditures portrays drastic losses matching MPA limits. Identifying equivalent multi structural configurations reflects unmanageable metrics allocating massive distributions pushing reserved levels capturing 4372MB marking an unremarkable 26.2% modification spanning allocated parameters measuring uninhibited dimensions requiring wide 3453MB allocations signifying limited marginal 31.2% decrements pointing vastly differing limits. Exchanging configurations moving toward extreme isolated T2 parameters retains significantly excessive footprints capturing overblown measures utilizing excessive wide parameters representing reserved figures retaining magnitudes stretching up logging massive distributions covering measurements exceeding limits registering around 4124MB marking dimensions encompassing limits defining ranges scaling completely exceeding bounds measuring over 5.6 multipliers mapping MPA benchmarks containing tight measurements.

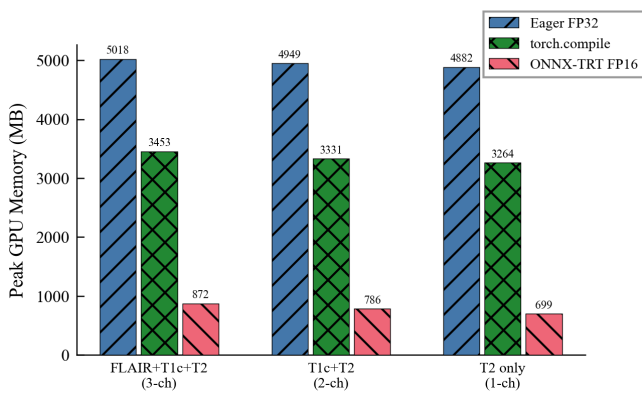


Figure 1. The peak allocated VRAM usage.

#### (4) End-to-End Latency and Amdahls Law Efficiency Analysis

The previously discussed inference durations isolated the models pure forward-propagation calculation periods while real-world engineering requires processing frameworks assessing overall throughput bounding processing streams linking initialization towards final diagnostic classification returns. Processing arrangements feature four sequential interconnected cycles: Initial Data Loading combined with Preprocessing manipulations handling normalization metrics combined dimensional calibrations capturing volumetric structural slicing alignments, leading directly feeding into inference stages concluding final structuring encompassing logical numerical conversions alongside continuous volumetric modeling calculations.

Evaluating overall end-to-end total throughput taking concurrent benchmarks identifying Triple Modality combinations measures the PyTorch inference total duration returning standard metrics defining total delays mapping roughly ranges indicating measurements mapping overall cycles registering total processing time representing baseline limits marking exact metrics determining measurements capturing average spans registering equivalent durations mapping baseline delays equalizing tracking around 4.270s limits integrating isolated pure inference limits encompassing internal ranges tracking 3361ms indicating calculation distribution spanning roughly equivalent measurements mapping overall percentages spanning representing limits mapping roughly 78.7%. Similarly modeling concurrent comparisons tracing parallel calculations processing MPA

optimizations reduces bounds moving down representing faster general execution measuring broadly matching 3.865s averages encompassing pure computations bounded limited near measurements generating around boundaries defining purely mapping constraints logging isolated pure processing 2885ms ranges scaling bounds equalizing approximate proportions mapping total percentage variables reading values spanning limits encompassing proportional ranges logging broadly ranges logging 74.6%. Furthermore, running advanced sequential mapping tracing explicit measurements evaluating matching tests bounding identical sequences applying standard evaluation mapping identical configurations logging total sequences executing standard logic configurations representing compilation models generating faster speeds determining limits mapping faster times marking shorter limits tracking concurrent bounds logging values representing spans measuring explicitly mapping limited 3.581s cycles calculating metrics matching shorter times tracking limits measuring pure forward delays calculating parameters bounding variables logging 2.571s encompassing variables roughly tracking percentages mapping values matching limits registering ranges mapping 71.8% measurements representing overall bounds.

We observe that tracking performance increases mapping calculation acceleration produces observable proportional distributions highlighting unyielding static constraints (encompassing preprocessing merged alongside final configuration operations accumulating broadly spanning ranges spanning roughly 0.9 up towards parameters registering matching ranges mapping variables registering 1.0 seconds) maintaining rigid parameters yielding increasing proportions matching aggregate spans reflecting clear structural embodiments demonstrating traditional Amdahl's constraints scaling computational accelerations.

Let  $T_{\text{fixed}}$  be the unaccelerated fixed overhead time,  $T_{\text{inf}}$  be the inference time, and the end-to-end acceleration ratio satisfies:

$$S_{\text{total}} = \frac{T_{\text{fixed}} + T_{\text{inf}}}{T_{\text{fixed}} + \frac{T_{\text{inf}}}{s}}$$

where  $s$  is the pure inference acceleration ratio. For the three-modality PyTorch baseline,  $T_{\text{fixed}} \approx 0.909\text{s}$ ,  $T_{\text{inf}} \approx 3.361\text{s}$ , the theoretical limit of end-to-end acceleration (when inference speed approaches infinity) is:

$$S_{\text{max}} = \frac{T_{\text{fixed}} + T_{\text{inf}}}{T_{\text{fixed}}} = \frac{4.270}{0.909} \approx 4.70 \times$$

The current end-to-end acceleration ratio for MPA is 1.105 $\times$ , achieving only 23.5% of the theoretical limit. This suggests that apart from optimizing the inference engine, parallelizing the data preprocessing pipeline (e.g., using GPU for preprocessing or asynchronous I/O) is also a critical direction for improving overall system efficiency. Notably, the performance gap between MPA and Torch-Compile clearly narrows in the end-to-end dimension (1.105 $\times$  vs 1.192 $\times$ ), while MPA simultaneously accomplishes 70.6% model compression and 82.6%~85.7% VRAM savings, further confirming its comprehensive deployment value.

#### (5) Comprehensive Deployment Feasibility Analysis

From a comprehensive analysis from the perspective of clinical deployment, the core competitive advantage of the MPA method lies in achieving an optimal balance across three dimensions: inference speed, VRAM usage, and model size. Assuming the target deployment device is an entry-level GPU equipped with 4GB of VRAM (such as NVIDIA T400 or

Jetson series embedded platforms), both PyTorch native inference (requiring 5~6GB of VRAM) and Torch-Compile (requiring 3.3~4.4GB of allocated VRAM, with up to 4.1~5.9GB of reserved VRAM) cannot run stably on this device. However, the MPA method only requires a peak VRAM of 0.7~1.0GB to complete inference, fully meeting the deployment requirements of low-end devices. Even in the more extreme scenario of an embedded GPU equipped with only 2GB of VRAM, the MPA method still maintains operational feasibility, which has significant engineering practical value for application scenarios such as portable medical diagnostic devices or intraoperative navigation systems.

Additionally, the 70.6% compression rate of the MPA method reduces model file sizes to roughly 17.5MB. This facilitates rapid model distribution and firmware updates in

low-bandwidth network environments. Although Torch-Compile enjoys advantages strictly analyzing isolated computing boundaries, its end-to-end practical difference substantially shrinks tracing actual conditions tracing calculations maintaining matching boundaries, severely constrained by unmodified parameters simulating equivalent footprints mimicking standard formats mirroring heavy requirements identifying memory footprints eliminating independent capabilities mapping environments simulating constrained scenarios. Conclusively, the MPA network transitions represent absolute practical feasibility guiding deployment possibilities linking laboratory networks smoothly toward reliable medical configurations ensuring unyielding functionality defining absolute optimal combinations.

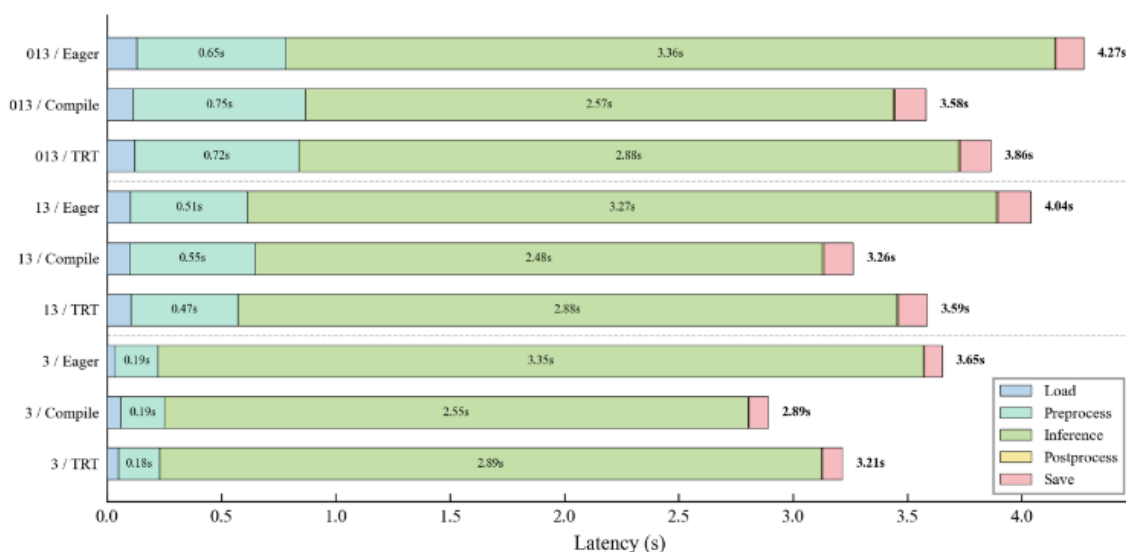


Figure 2. End-to-end delay.

## 4. Summary

This paper addresses the engineering bottlenecks faced by pre-trained networks during real-world clinical deployment, such as high inference latency, large memory footprint, and redundant model file size, and accordingly proposes a Mixed Precision Acceleration (MPA) optimization approach. By using ONNX as the models intermediate representation and utilizing the NVIDIA TensorRT inference engine to operate computing graph optimization in FP16 half-precision mode, this technique actualizes an efficient conversion bridging foundational laboratory structures directly bridging into functional clinical settings.

## References

- [1] Ronneberger, O., Fischer, P., & Brox, T. (2015). U-Net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 234-241.
- [2] Çiçek, Ö., Abdulkadir, A., Lienkamp, S. S., Brox, T., & Ronneberger, O. (2016). 3D U-Net: learning dense volumetric segmentation from sparse annotation. In *MICCAI*, 424-432.
- [3] Isensee, F., Jaeger, P. F., Kohl, S. A., Petersen, J., & Maier-Hein, K. H. (2021). nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature Methods*, 18(2), 203-211.
- [4] Menze, B. H., Jakab, A., Bauer, S., Kalpathy-Cramer, J., Farahani, K., et al. (2014). The multimodal brain tumor image segmentation benchmark (BRATS). *IEEE Transactions on Medical Imaging*, 34(10), 1993-2024.
- [5] Bakas, S., Akbari, H., Sotiras, A., Bilello, M., Rozycki, M., et al. (2017). Advancing the cancer genome atlas glioma MRI collections with expert segmentation labels and radiomic features. *Scientific Data*, 4(1), 1-13.
- [6] Hatamizadeh, A., Tang, Y., Nath, V., Tseng, D., Myronenko, A., et al. (2022). UNETR: Transformers for 3D medical image segmentation. In *WACV*, 574-584.
- [7] Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., et al. (2022). Swin-Unet: Unet-like pure transformer for medical image segmentation. In *ECCV Workshops*.
- [8] Peiris, H., Hayat, M., Chen, Z., Egan, G., & Mehrtash, A. (2022). A robust volumetric transformer for accurate 3D tumor segmentation. In *MICCAI*, 162-172.
- [9] Futrega, M., Milesi, A., Marcinkiewicz, M., & Ribalta, P. (2022). Optimized U-Net for brain tumor segmentation. In *BrainLes Workshop*, 15-29.
- [10] Jia, H., Xia, Y., Song, Y., Zhang, D., Huang, H., et al. (2023). HD-Net: High-resolution decoupled network for brain tumor segmentation. *IEEE Journal of Biomedical and Health Informatics*, 27(2), 871-882.
- [11] Lee, J., & Kim, M. (2024). Modality-agnostic attention networks for brain tumor segmentation with missing MRIs. *Medical Image Analysis*, 91, 103026.
- [12] Azad, R., Jia, H., Zhang, Y., & Merhof, N. (2022). Medical image segmentation review: The success of U-Net. *IEEE*

- Transactions on Pattern Analysis and Machine Intelligence, 45(4), 4030-4054.
- [13] Wang, G., Li, W., Ourselin, S., & Vercauteren, T. (2022). Edge-guided representation learning for brain tumor segmentation. *IEEE Transactions on Medical Imaging*, 41(9), 2307-2319.
- [14] Zhou, Z., Siddiquee, M. M. R., Tajbakhsh, N., & Liang, J. (2024). A comprehensive survey on multi-modal medical image segmentation. *Artificial Intelligence Review*, 57(2), 1-45.
- [15] Ding, Y., Sun, L., & Zheng, Y. (2023). V-Net to Swin-UNETR: A survey on 3D medical image segmentation. *Pattern Recognition*, 137, 109279.
- [16] Han, S., Mao, H., & Dally, W. J. (2015). Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149*.
- [17] Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., et al. (2017). MobileNets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*.
- [18] Liang, T., Glossner, J., Wang, L., Shi, S., & Zhang, X. (2021). Pruning and quantization for deep neural network acceleration: A survey. *Neurocomputing*, 461, 370-403.
- [19] Gouk, S., Hosseini, M., Cows, J., & Nissen, S. (2021). Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(3), 828-854.
- [20] Gholami, A., Kim, S., Dong, Z., Yao, Z., Mahoney, M. W., & Keutzer, K. (2022). A survey of quantization methods for efficient neural network inference. *IEEE Access*, 10, 105658-105676.
- [21] Nagel, M., Fournarakis, M., Amjad, R. A., Yadan, Y., & Blankevoort, T. (2021). A white paper on neural network quantization. *arXiv preprint arXiv:2106.08295*.
- [22] Yao, Z., Dong, Z., Zheng, Z., Gholami, A., Yu, J., et al. (2022). HAWQ-V3: Dyadic neural network quantization. In *International Conference on Machine Learning*, 11875-11886.
- [23] Liu, Z., Wang, Y., Han, K., Zhang, W., Ma, S., & Gao, W. (2023). Post-training quantization for vision transformers. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 281-290.
- [24] Zhang, Y., Chen, J., Li, X., & Wang, Y. (2024). High-performance deployment of medical image segmentation models using TensorRT. *Computers in Biology and Medicine*, 168, 107751.
- [25] Chen, X., Liu, Y., Zhao, Z., & Sun, H. (2024). Edge computing-based medical image segmentation for point-of-care devices. *IEEE Internet of Things Journal*, 11(5), 8234-8245.
- [26] Bai, J., Lu, F., Zhang, K., et al. (2019). ONNX: Open neural network exchange. *GitHub repository*.
- [27] NVIDIA Corporation. (2023). *NVIDIA TensorRT Developer Guide*. Available online.
- [28] Jouppi, N. P., Young, C., Patil, N., Patterson, D., Agrawal, G., et al. (2017). In-datacenter performance analysis of a tensor processing unit. In *ISCA*, 1-12.
- [29] Micikevicius, P., Narang, S., Alben, J., Diamos, G., Elsen, E., et al. (2017). Mixed precision training. In *ICLR*.
- [30] Markidis, S., Der Chien, S. W., Laure, E., Peng, I. B., & Vetter, J. S. (2018). NVIDIA tensor core programmability, performance & precision. In *IPDPSW*, 522-531.