





Multimodal Diverse Granularity Fusion Network based on US and CT Images for Lymph Node Metastasis Prediction of Thyroid Carcinoma

Guojun Li^{1, 2*}, Jincao Yao^{3*}, Chanjuan Peng^{3*}, Yinjie Hu²,
Shanshan Zhao⁴, Xuhan Feng², Jianfeng Yang⁴,
Dong Xu³ , Xiaolin Li^{1, 2} , Chulin Sha² , Min He² 

¹Academy of Medical Engineering and Translational Medicine, Tianjin University

²Hangzhou Institute of Medicine, Chinese Academy of Sciences

³Zhejiang Cancer Hospital

⁴Shaoying People's Hospital

Abstract

Accurately predicting the risk of cervical lymph node metastasis (LNM) is crucial for surgical decision-making in thyroid cancer patients, and the difficulty in it often leads to over-treatment. Ultrasound (US) and computed tomography (CT) are two primary non-invasive methods applied in clinical practice, but both contain limitations and provide unsatisfactory results. To address this, we developed a robust and explainable multimodal deep-learning model by integrating the above two examinations. Using 3522 US and 7649 CT images from 1138 patients with biopsy-confirmed LNM status, we showed that multimodal methods outperformed unimodal counterparts at both central and lateral cervical sites. By incorporating a diverse granularity fusion module, we further enhanced the area under the curve (AUC) to 0.875 and 0.859 at central and lateral cervical sites respectively. This performance was also validated in an external cohort. Additionally, we quantified the modality-specific contributions for each nodule and systematically evaluated the applicability across various clinical characteristics, aiding in identifying individuals who can benefit most from the multimodal method.

1 The global incidence of thyroid cancer has surged
2 over the past 30 years[1], reaching over 586,000 new
3 cases in 2020[2]. Despite its generally indolent nature,
4 thyroid cancer leads to cervical lymph node
5 metastasis (LNM) in up to 50% of patients[3]. Cancer
6 cells typically initially metastasize to the central
7 lymph nodes and subsequently spread to the lateral
8 cervical site, increasing the risk of recurrence and
9 poor prognosis[4]. Consequently, LNM status significantly
10 influences the surgical approach for thyroid cancer patients.
11 Therapeutic lymph node dissection (LND) of central and lateral
12 cervical compartments is normally recommended for individuals
13 with central and/or lateral cervical LNM[5]. While for patients
14 without LNM, although central LND remains controversial,
15 prophylactic lateral cervical LND is not advised[5, 6].
16 However, the current non-invasive diagnostic accuracy of LNM
17 is insufficient to guide surgical decisions. For the central site,
18 the primary imaging methods, including Ultrasound (US) and
19 computed tomography (CT), provide average sensitivities of
20 only 0.28 and 0.39[7], respectively. This leads to a prevalent
21 tendency for overtreatment to prevent missed LNM detection
22 and results in potential complications such as recurrent laryngeal
23 nerve

injury and hypoparathyroidism. Therefore, there is a
pressing need to improve the accuracy of LNM risk
assessment to assist surgical management.

In recent years, the introduction of artificial intelligence
methods has significantly improved the performance of LNM
prediction. Several studies utilizing US images have employed
various machine learning methods, such as gradient boosting,
random forests, neural networks, etc., achieving AUCs in the
range of 0.700 to 0.772[8, 9, 10] for predicting central site
LNM. Other studies focusing on extracting high-dimensional
radiomic features or employing deep learning methods to
predict LNM status have achieved AUCs spanning from 0.78 to
0.90[11, 12, 13] for the central site and 0.62[14] for the
lateral cervical site. Similarly, in the case of CT images,
methods based on radiomic features extracted from thyroid
nodules have demonstrated predictive capabilities for central
site LNM at AUCs ranging from 0.710 to 0.770[15, 16].

However, it's crucial to acknowledge that both US and CT
modalities have limitations owing to their examination
techniques. Though US images provide high-resolution
visuals of thyroid nodules' interior and boundary
characteristics, their limited field of view poses challenges
in assessing the spatial relationships between nodules and
surrounding tissues. Conversely, CT images offer essential
relative

* These authors contributed equally to this work.

† Corresponding authors: xudong@zjcc.org.cn, xiaolin@icm.cas.ac.cn, shachulin@gmail.com, hemin@him.cas.cn

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

54 position information about the thyroid, lymph nodes,
55 and surrounding tissues, albeit at a lower resolution
56 compared to US images. Relying solely on unimodal
57 methods restricts the predictive capabilities of the
58 model. Considering that US and CT provide comple-
59 mentary information and are widely utilized in
60 thyroid cancer diagnosis, there's a great potential to
61 improve the performance by integrating US and CT
62 images through a multimodal approach for predict-
63 ing LNM status. For instance, Zhao et al. developed
64 a multivariate logistic regression multimodal model
65 for predicting central LNM status by incorporating
66 clinical factors, US-derived diagnostic features, and
67 CT measurements, achieving an AUC of 0.827[17].
68 Nevertheless, the study did not directly compare the
69 multimodal method's performance with unimodal
70 methods, besides, it utilized a simplified model that
71 overlooked the interaction between the two modal-
72 ities, leaving the full potential of multimodal fusion
73 approach unexplored.

74 Leveraging deep learning methods for integrating
75 multimodal medical data has emerged as a promi-
76 nent approach to enhance our understanding of com-
77 plex diseases[18, 19, 20], with promises in tailoring
78 personalized diagnosis, prognosis, treatment, and
79 care[21, 22, 23, 24]. The central premise of multi-
80 modal data integration is that diverse data sources
81 complement each other, augmenting information be-
82 yond any individual modality. However, significant
83 challenges persist, such as data scarcity, sparsity, and
84 inter-modality complexity, limiting the full exploita-
85 tion of data integration benefits. Recent advance-
86 ments in deep learning methods within this domain
87 primarily focus on representation learning and fusion
88 techniques[25, 26], which include extracting mean-
89 ingful representations with unlabeled data[27, 28]
90 and employing attention-based approaches to allow
91 more sophisticated fusion of cross-modality represen-
92 tations[29, 30, 31, 32]. While these strategies exhibit
93 improvements in model performance, their use in the
94 biomedical field still requires broader testing across
95 diverse scenarios and adaptation to specific tasks
96 through dedicated study designs[33].

97 In this study, we aim to improve the predictive
98 accuracy of LNM status for thyroid cancer patients
99 by developing a multimodal method incorporating
100 US and CT images. We curated a paired multi-
101 modal dataset consisting of 3522 US and 7649 CT
102 images from 1138 patients with biopsy-confirmed
103 LNM status at both central and lateral cervical sites
104 (Fig. 1). To comprehensively integrate the consis-
105 tent and distinct information of both modalities, we
106 first employed a multi-task network scheme to en-
107 hance modal-specific feature learning (Fig. 2), which
108 achieved superior performance compared to cur-

rently commonly used methods on unimodal mod- 109
els. Next, we demonstrated that, even with a basic 110
feature fusion strategy, multimodal models consis- 111
tently outperform their unimodal counterparts at 112
both sites. Furthermore, we designed a diverse gran- 113
ularity fusion module, which learns the attention at 114
three granular levels from fine to coarse: dimension 115
level, modality level, and nodule level (Fig. 2). With 116
the incorporation of this module, our multimodal 117
model achieved AUCs of 0.875 and 0.859 at the cen- 118
tral and lateral cervical sites respectively. Compared 119
to unimodal methods of US and CT, the multimodal 120
AUC improved by 5.5% and 10.1% respectively, at 121
the central compartment, and by 7.4% and 8.1% re- 122
spectively, at the lateral cervical site. When evaluated 123
on an external validation set, our proposed model 124
demonstrated an AUC of 0.903 at the central site, 125
which robustly confirmed the generalizability of the 126
multimodal model. In addition, we comprehensively 127
evaluate the applicability of each modality on nod- 128
ules with various characteristics to identify patients 129
who can best benefit from the multimodal method 130
(Fig. 1), which could significantly improve the clin- 131
ical utility of multimodal models. In summary, we 132
presented a promising approach to mitigate the issue 133
of overtreatment in thyroid cancer. Our multimodal 134
AI system exhibits strong performance, high gen- 135
eralizability, and substantial clinical utility, offering 136
significant potential for enhancing the diagnosis and 137
treatment of thyroid cancer. 138

139 Results

140 Patient Cohort

141 This study incorporated two datasets: a main cohort
142 and an external cohort. The main cohort comprised
143 patients who underwent thyroid examinations at Zhe-
144 jiang Cancer Hospital from August 2018 to February
145 2021. To reflect the real clinical diagnostic conditions,
146 only necessary data quality control was performed,
147 with specific details outlined in the supplementary
148 material. After quality control, the main cohort con-
149 sists of 1138 patients with a total of 1285 thyroid nod-
150 ules. The external cohort, obtained from Shaoxing
151 People's Hospital in Zhejiang Province, also under-
152 went the same quality control process, comprising
153 60 patients with 60 thyroid nodules. Both cohorts
154 included samples with matched US and CT data,
155 featuring multiple images of thyroid nodules, along
156 with their corresponding LNM status at the central
157 and lateral cervical sites.

158 The models were evaluated under an eight-fold
159 cross-validation setting, and various metrics were
160 employed to assess their performance. These met-

Table 1: Performance of LNM status prediction using unimodal networks

Modalities	LNM location	Models	ACC	AUC	SENS	SPEC	PREC	F1 score
US	Central site	Yu et al [12]	0.739	0.792	0.787	0.689	0.719	0.750
		ResNet	0.760	0.812	0.823	0.699	0.736	0.774
		Proposed	0.782	0.820	0.836	0.728	0.754	0.792
	Lateral cervical site	Yu et al [12]	0.729	0.754	0.729	0.731	0.736	0.728
		ResNet	0.747	0.725	0.826	0.665	0.713	0.763
		Proposed	0.767	0.785	0.843	0.693	0.737	0.784
CT	Central site	ResNet	0.726	0.759	0.743	0.707	0.720	0.730
		Proposed	0.745	0.774	0.737	0.755	0.751	0.743
	Lateral cervical site	ResNet	0.746	0.758	0.812	0.679	0.718	0.760
		Proposed	0.764	0.778	0.725	0.797	0.796	0.751

rics encompassed accuracy (ACC), the area under the ROC curve (AUC), sensitivity (SENS), specificity (SPEC), precision (PREC), and the F1 score. We reported the mean metrics calculated from the eight-fold cross-validation process for a comprehensive evaluation. It is worth mentioning that samples from different folds were divided based on the individual thyroid nodules, and nodules from the same patient were consistently present within the same fold. In addition, an undersampling strategy was applied in this study to maintain a balance between positive and negative categories.

Enhance modal-specific feature learning by employing multi-task models of each modality

We start from enhancing modal-specific feature extraction to make the best use of each modality and evaluate the feature capability in predicting LNM status of each modality. US provides clear visualization of thyroid nodule attributes such as boundary, shape, and internal structure (composition, calcification, echo characteristics). Meanwhile, CT images encompass both thyroid nodules and the surrounding anatomical context, offering insights into their relationships. Therefore, we employ a multi-task learning approach for each modality (as illustrated in Fig. 2). Specifically, besides the LNM prediction task, we introduce two auxiliary tasks for US: a nodule mask segmentation task to guide the model to focus on the internal structural features, and a nodule boundary segmentation task to emphasize the boundary and shape of nodules. Likewise, for CT, we introduce a nodule mask task and a tissue boundary segmentation task to guide the model to distinguish nodules and surrounding tissue regions, respectively. We chose ResNet[34] as the backbone to build multi-task models for each modality, due to its simple structure, high popularity, and excellent performance (Methods). For each unimodal network,

we trained it to complete the auxiliary segmentation tasks in the first 100 epochs and added the additional LNM prediction task in the following 200 epochs. We compared our multi-task models for each modality with ResNet models directly predicting LNM, and for the US unimodal model, we also re-implemented the network developed by Yu et al[12]. It shows that our multi-task models for both modalities consistently outperform their counterparts at both central and cervical lateral sites, with an obvious improvement of ACC and AUC (Table 1).

When comparing the unimodal performance of US and CT, we have some interesting observations. First, at the central site, the US models generally outperform CT models, whereas there is no consistent winner between US and CT models at the lateral cervical site. Moreover, at both the central and lateral cervical sites, US models consistently exhibit higher sensitivity, meanwhile, CT models consistently demonstrate higher specificity. These results suggest that US is more sensitive but less specific, while CT is the opposite, highlighting the complementary information provided by these two modalities.

Basic multimodal fusion methods outperform either unimodal model

Based on the multi-task unimodal models, we further evaluate the efficacy of integrating US and CT for predicting the risk of LNM. We first examine the multimodal performance using three basic fusion methods: concatenation, element-wise sum, and element-wise multiplication, to fuse the unimodal features extracted from US and CT and re-train the multimodal network in an end-to-end manner. The results clearly show that even with basic fusion methods, multimodal models significantly improve performance.

For the central site prediction, the average multimodal AUC improved by 2.8% and 7.4% compared to US and CT unimodal respectively. Likewise,

Table 2: Performance comparison of unimodal and multimodal approaches using basic fusion methods

LNM location	Modality	Multimodal fusion operation	ACC	AUC	SENS	SPEC	PREC	F1 score
Central site	US	Concat	0.782	0.820	0.836	0.728	0.754	0.792
			0.745	0.774	0.737	0.755	0.751	0.743
	Multimodal	Concat	0.810	0.853	0.862	0.759	0.782	0.819
		Sum	0.806	0.850	0.855	0.757	0.779	0.814
		Product	0.812	0.840	0.861	0.763	0.790	0.820
		Average results	0.809	0.848	0.859	0.760	0.784	0.818
Lateral cervical site	US	Concat	0.767	0.785	0.843	0.693	0.737	0.784
			0.764	0.778	0.725	0.797	0.796	0.751
	Multimodal	Concat	0.804	0.825	0.860	0.744	0.779	0.815
		Sum	0.808	0.854	0.878	0.742	0.771	0.819
		Product	0.811	0.835	0.867	0.755	0.784	0.821
		Average results	0.808	0.838	0.868	0.747	0.778	0.818

for the lateral cervical site prediction, the average multimodal AUC outperforms the US and CT unimodal models by 5.3% and 6.0% respectively (Table 2). These results affirm the hypothesis that US and CT modalities comprise complementary information, and their integration can improve the performance of LNM status prediction.

Further improve multimodal performance by incorporating a diverse granularity feature fusion module

Multimodal fusion using basic methods can combine US and CT information and improve LNM prediction but is not able to fully consider the interaction between these two modalities. Recent progress based on the attention mechanism has shown superiority in multimodal fusion. In our study, we adopted the attention mechanism simultaneously on three granular levels to fully incorporate the information useful for LNM prediction, which are feature dimensions level (minimum granularity), modalities level (medium granularity), and nodules level (maximum granularity). In specific, these include dynamically adjusting the attention weights of different feature dimensions to balance the common and specific features of the two modalities, adapting the modality-specific attention to learn the respective advantages for different nodules, plus flexibly aggregate the features of other nodules based on nodule-level attention to refine the prediction, considering nodules with the same LNM status should exhibit greater feature similarity. We refer to our modality fusion methods as the ‘diverse granularity fusion’ network (DGFNet, as illustrated in Fig. 2, detail see Methods). Equipped with the DGF module, our model demonstrates exceptional predictive capabilities with AUCs of 0.875 and 0.859 at the central and lateral cervical sites respectively

(Table 3), indicating its remarkable performance in predicting the risk of LNM. Particularly, the multimodal AUC exhibited a significant improvement of 5.5% and 10.1% compared to the US and CT unimodal models at the central site, respectively. And a substantial enhancement of 7.4% and 8.1% respectively at the lateral cervical site. These results further underscore the efficacy of integrating US and CT in predicting LNM in thyroid cancer. Furthermore, in comparison to the basic fusion methods, our DGFNet model achieves superior performance in nearly all metrics, providing comprehensive evidence for the effectiveness of fusing multimodal features at different granularities.

DGFNet demonstrates exceptional generalization abilities

Recognizing the importance of the generalizability of multimodal networks in clinical applications, we evaluated the efficacy of our DGFNet model using an external test dataset. The primary cohorts were partitioned into training and validation sets, and the model with the highest accuracy on the validation set was selected to predict the LNM status of patients on the external cohort. Owing to constraints related to data availability, our external evaluation is only performed at the central site, the results are presented in Table 4. Overall, our DGFNet model performed well on the external dataset, with an accuracy of 0.817 and an AUC of 0.903, showing similar performance compared to the internal accuracy of 0.844 and an AUC of 0.898. This consistency underscores the strong robustness and external generalizability of our model.

Table 3: Performance comparison of multimodal methods using different fusion techniques

LNM location	Multimodal fusion methods	ACC	AUC	SENS	SPEC	PREC	F1 score
Central site	Basic fusion	0.809	0.848	0.859	0.760	0.784	0.818
	Diverse granularity fusion	0.826	0.875	0.848	0.803	0.813	0.830
Lateral cervical site	Basic fusion	0.808	0.838	0.868	0.747	0.778	0.818
	Diverse granularity fusion	0.838	0.859	0.862	0.814	0.835	0.842

Table 4: Performance of LNM status prediction on internal and external validation set

Dataset	ACC	AUC	SENS	SPEC	PREC	F1 score
Internal validation set	0.844	0.898	0.854	0.833	0.837	0.845
External validation set	0.817	0.903	0.909	0.763	0.690	0.784

DGFNet dynamically adjusts the contribution of US and CT in predicting the LNM status prediction at different sites

We next seek to delineate the contribution of each modality in the DGFNet model on every nodule. We analyze by quantifying the contributions of US and CT within the DGFNet model using the integrated gradients[35] and comparing them to their unimodal counterparts. The results are presented in Fig. 3, where a larger feature attribution value corresponds to a greater contribution to the correct prediction in DGFNet model, and the red or green denotes correct or wrong prediction respectively.

The result shows that, at both the central site (Fig. 3a) and lateral cervical site (Fig. 3b), there is a notable number of cases where the DGFNet can change the unimodality to make positive contributions even when it fails to give correct prediction in unimodal models (attribution greater than 0 but red color) and lead to a correct prediction in this multimodal approach. In addition, when looking at the central and lateral cervical sites separately, we find that, across all samples, 64.9% of nodules exhibit higher US attribution over CT attribution at the central site, while 55.2% of the nodules show higher US attribution over CT attributions at the lateral cervical site. These findings agree with our prior observations on unimodal LNM prediction performance, highlighting a more prominent role for US at the central site, whereas both US and CT show comparable importance at the lateral cervical site. In addition, this underscores that the DGFNet model can dynamically adjust the weights of the two modalities based on nodule characteristics, effectively leveraging the strengths of both modalities.

DGFNet enhances model attention on the nodular region in US and CT images

To further investigate how our DGFNet model improves the LNM prediction performance, we generated saliency maps for both US and CT images in the multimodal network and compared them with their unimodal counterparts. The results clearly show that, for both US and CT images, the DGFNet model significantly increases the attention towards the region of interest compared to the unimodal models. Specifically, within US images, the multimodal model focuses more intensely on the nodules' peripheral and inner hypoechoic region (Fig. 4a), whereas in CT images, it narrows its focus to the nodules and their immediate surrounding tissues (Fig. 4b), all of which represent crucial regions providing key information for LNM prediction. This directly proves the superiority of DGFNet in grasping meaningful medical information over unimodal methods.

Identify patients who can best benefit from multimodal integration

The multimodal approach can effectively improve the LNM prediction, however, it is often unfeasible to examine all patients by both modalities in real clinical settings. Hence, to make our DGFNet more applicable and useful for clinicians, we further seek to identify patients who can best benefit from the multimodal approach. Given that the US examination is cheaper and more commonly used, we analyzed by identifying cases for whom adding CT as a supplementary modality would be advantageous. We evaluated four well-established sonographic characteristics of the thyroid nodule during US diagnosis including maximum diameter, margin characteristics, aspect ratio, plus location in the thyroid for central cite nodules, and categorized the nodules based on the measurements. We then compared the prediction performance in each category between the DGFNet

382 and the unimodal model of US and CT respectively.
383 (Results on more characteristics are illustrated in Sup-
384 plementary Material)

385 The analysis shows that the DGFNet model
386 achieves particularly high performance in specific
387 circumstances. This includes nodules with maxi-
388 mal diameters between 20mm and 36mm, as well as
389 more extreme cases less than 12mm or larger than
390 60mm(Fig. 5a). Additionally, DGFNet excels in cases
391 exhibiting non-smooth borders (Fig. 5b), aspect ratios
392 surpassing 1 (Fig. 5c), and nodules situated within
393 the thyroid isthmus (Fig. 5d). Similar findings are
394 observed in the analysis conducted at the lateral cer-
395 vical site (Supplementary Material). Therefore, the
396 DGFNet model is potentially particularly beneficial
397 and practical for patients with the above nodule char-
398 acteristics.

399 Discussion

400 Patients often undergo multiple types of examina-
401 tions in the diagnostic process, and the effective in-
402 tegration of multimodal information can greatly im-
403 prove diagnosis accuracy. Recent advancements in
404 artificial intelligence techniques have facilitated the
405 progress of deep-learning-based multimodal integra-
406 tion methods, which have emerged as a trend in
407 cancer diagnosis in recent years. In this study, we
408 pioneered the development of a multimodal deep
409 learning approach that effectively integrates US and
410 CT modalities to successfully enhance the accuracy
411 of LNM prediction and further demonstrate its gen-
412 eralizability in an external dataset. Moreover, by
413 conducting a series of comprehensive interpretability
414 analyses, we quantified the modality-specific con-
415 tribution across nodules in various situations, and
416 investigated the attention heatmap of US and CT im-
417 ages within the model, which not only shed light
418 on the reasons for the improved performance of the
419 multimodal model, but also improve the model's
420 applicability in clinical settings, and opens a new
421 avenue for mitigating the problem of overtreating
422 thyroid cancer.

423 The effective integration of multimodal data often
424 relies on a deep understanding of the domain knowl-
425 edge involved with specific medical tasks. In our
426 study, a close collaboration between AI scientists and
427 clinicians allowed us to leverage our collective exper-
428 tise in deep learning models, thyroid cancer, US and
429 CT images. This enabled us to strategically employ
430 multi-task learning techniques, facilitating the identi-
431 fication of critical regions and extraction of essential
432 LNM-related features from both US and CT images.
433 Moreover, we introduced a novel diverse granularity
434 fusion network (DGFNet) that learns the attention

435 from three different levels, which excels in not only
436 effectively integrating shared and specific features
437 from multimodal data but also dynamically adjust-
438 ing the weights of each modality's data for different
439 nodules. This approach demonstrates the potential
440 to optimize the utility of both US and CT images and
441 aggregate information from similar nodules, thereby
442 enhancing the model's overall performance and ro-
443 bustness.

444 Besides the excellent performance of our developed
445 DGFNet model, our study has yielded valuable clinical
446 insights through the multimodal approach. First,
447 it shows that unimodal methods based on US appear
448 to be more sensitive but less specific, while CT-based
449 unimodal methods are the other way around. Second,
450 it shows that the US modality generally plays a more
451 significant role than CT at the central site, whereas
452 there is no obvious difference between US and CT at
453 the lateral cervical site. Furthermore, by quantifying
454 the performance of the unimodal and multimodal
455 models for nodules within different diagnosis char-
456 acteristics categories, we could pinpoint patients with
457 certain nodule characteristics who can potentially
458 best benefit from the multimodal approach. These
459 analyses offer valuable insights for accurately iden-
460 tifying the appropriate patient population for mul-
461 timodal diagnostic approaches in clinical practice
462 and guiding patients in selecting the most suitable
463 examination method.

464 In conclusion, through a close collaboration be-
465 tween AI scientists and clinicians, this study suc-
466 cessfully develops a multimodal approach aimed at
467 improving the LNM prediction for thyroid patients.
468 It paves the way for addressing the issue of overtreat-
469 ment in thyroid cancer and provides new insights
470 in the integration of multimodal data for precise di-
471 agnosis, representing an excellent scientific research
472 example originating from clinical practice and di-
473 rectly addressing clinical necessities.

474 Methods

475 Patient Cohort

476 There are two cohorts included in this study. The
477 main cohort was obtained from Zhejiang Cancer Hos-
478 pital in Zhejiang Province, China, consisting of 1360
479 patients. After the data screening process, a total of
480 1138 patients with 1285 nodules were retained for
481 analysis. The main cohort was utilized for the model
482 establishment and internal performance evaluation.
483 The second cohort, referred to as the external co-
484 hort, was sourced from Shaoxing People's Hospital
485 in Zhejiang Province, China. Initially, this cohort in-
486 cluded 126 patients, and after the data screening pro-

487 cess, 60 patients with complete data were included
488 for evaluation of model generalization (The patient
489 enrollment process is illustrated in Supplementary
490 Material). Ethical approval for the study was ob-
491 tained from the ethics committees of both hospitals
492 and verbal informed consent was obtained from all
493 participating patients.

494 The inclusion criteria for this study encompassed
495 the following: (1) patients with thyroid nodules, (2)
496 patients who underwent cervical US and CT exami-
497 nations, and (3) patients with confirmed pathological
498 status of cervical LNM. Exclusion criteria consisted
499 of the following: (1) missing US or CT data, (2) the
500 presence of measurement lines in the US images, and
501 (3) patients with multiple malignant thyroid nod-
502 ules and metastatic cervical lymph nodes. As all the
503 thyroid nodules had the potential to metastasize, it
504 was impossible to determine which specific nodule
505 had metastasized to the cervical lymph nodes. After
506 the data screening process, the number of nodules
507 with and without metastasis in the central and lat-
508 eral cervical sites for both cohorts is presented in the
509 Supplementary Material.

510 Multiple US images, including transverse and lon-
511 gitudinal sections, as well as multiple CT images
512 from different slices, were available for most nodules
513 in the dataset. During each epoch of the training
514 process, one random US image and one random CT
515 image were paired together to form an image pair.
516 During the evaluation process, the US and CT images
517 with the largest nodal area were selected from the
518 multiple available images to form an image pair for
519 analysis.

520 Data pre-processing

521 **Region of Interest Extraction.** The methods used
522 for extracting the region of interest in both US and
523 CT images are similar and described as follows: 1)
524 We first performed a dilation operation on the mask
525 of thyroid nodules annotated by clinicians, using a
526 3x3 dilation kernel. The iteration steps were set to
527 40 and 25 for US and CT images, respectively. 2) We
528 determined the horizontal bounding rectangle of the
529 dilated region, with the height, width, and center
530 coordinates of the rectangle denoted as h , w , and
531 $(x_{\text{center}}, y_{\text{center}})$, respectively. 3) Using $(x_{\text{center}}, y_{\text{center}})$
532 as the center and the larger value of h and w as the
533 side length, we obtained the external square of the
534 thyroid nodule. 4) The original US and CT images
535 were then cropped to reserve the region within this
536 square. If the square area exceeds the image bound-
537 ary, the images are padded with zeros to fill the
538 exceeding part.

539 **Image Augmentation.** The cropped US and CT
540 images were resized to 288x288 pixels and 96x96 pix-

els, respectively. To enhance the diversity of the data,
we applied additional data augmentation techniques
to both modalities. These techniques included rota-
tion, horizontal flip, cropping and scaling, brightnes-
s-contrast transformation, and elastic transformation.
For rotation, the angle of rotation ranged from -15°
to 15° . Random cropping occurred with the cropped
area set to be between 90% and 100% of the original
size. The probability of applying these transforma-
tions was set to 0.5, ensuring a balanced augmenta-
tion effect.

552 Convolutional neural network architecture

553 The architecture of the proposed model is depicted
554 in Fig. 2. The model is composed of three distinct
555 branches: the US branch, the CT branch, and the Mul-
556 timodal branch. Both the US and CT branches share
557 an identical structure, each comprising an encoder
558 and two decoders. The encoder adopts a pre-trained
559 ResNet[34] architecture, with ResNet34 and ResNet18
560 selected for the central and lateral cervical sites, re-
561 spectively. Regarding the US branch, the decoders
562 are trained to delineate the mask and boundary of
563 thyroid nodules, directing the model's attention to-
564 wards the internal and marginal regions of the nod-
565 ules, correspondingly. Conversely, the CT branch's
566 decoders focus on segmenting the mask of thyroid
567 nodules and the boundary of surrounding tissue, fa-
568 cilitating the model in comprehensively capturing
569 information about both the thyroid nodule and its
570 adjacent surroundings.

571 All the aforementioned decoders share the same
572 structure. Each decoder is constructed from 5 upsam-
573 ple blocks, with every block encompassing 2 layers.
574 In the initial layer of each block, the input feature
575 is upsampled using bilinear interpolation. Subse-
576 quently, the second layer comprises a convolutional
577 block, incorporating a convolutional layer featuring a
578 kernel size of 3×3 , followed by batch normalization,
579 relu activation, and a dropout layer. Notably, to en-
580 hance segmentation performance, short connections
581 interconnect the encoder and decoder components.

582 The US and CT encoders produce 512-dimensional
583 vectors through global average pooling. In unimodal
584 models, these vectors directly enter the classifier
585 for LNM prediction. In the multimodal model, the
586 unimodal vectors integrate within the multimodal
587 branch and then proceed to the classifier with the
588 same structure—a two-layer fully connected neural
589 network with 512 input nodes and a single output
590 node.

591 Diverse granularity fusion module

592 The diverse granularity fusion module comprises
593 three branches, as depicted in the supplementary
594 material. All branches are constructed using the
595 attention mechanism.

596 In the dimensional correlation branch, the US and
597 CT features undergo a preliminary transformation as
598 outlined below:

$$Q_{US}^D = f_{US} \times W_{Q-US}^D \quad (1)$$

$$V_{US}^D = f_{US} \times W_{V-US}^D \quad (2)$$

$$Q_{CT}^D = f_{CT} \times W_{Q-CT}^D \quad (3)$$

$$V_{CT}^D = f_{CT} \times W_{V-CT}^D \quad (4)$$

599 Here, f_{US} and f_{CT} are the unimodal features of US
600 and CT, respectively, and W_{Q-US}^D , W_{V-US}^D , W_{Q-CT}^D ,
601 and W_{V-CT}^D are trainable parameters. The product
602 of Q_{US}^D and K_{CT}^D , followed by the application of the
603 softmax function, results in the attention matrix A^D ,
604 which captures the interplay between various feature
605 dimensions of the US and CT modalities:

$$A^D = \text{softmax} \left(\frac{Q_{US}^D \times K_{CT}^D}{\sqrt{d_k}} \right) \quad (5)$$

606 Here, d_k is the dimension of K_{CT}^D . The derived
607 attention matrix is then utilized for the enhanced
608 multimodal features:

$$f_{US}^D = V_{US}^D + A^D \times V_{US}^D \quad (6)$$

$$f_{CT}^D = V_{CT}^D + A^{DT} \times V_{CT}^D \quad (7)$$

609 Ultimately, the enriched features are amalgamated
610 through concatenation along the dimension axis,
611 yielding the fused features:

$$f^D = \text{concat} \left(f_{US}^D, f_{CT}^D \right) \quad (8)$$

612 In the modal weights branch, the US and CT fea-
613 tures are first concatenated along the modal axis:

$$f_{US-CT}^N = \text{concat} \left(f_{US}, f_{CT} \right) \quad (9)$$

614 Then the Q^M , K^M , and V^M are generated respec-
615 tively:

$$Q^M = f_{US-CT}^M \times W_Q^M \quad (10)$$

$$K^M = f_{US-CT}^M \times W_K^M \quad (11)$$

$$V^M = f_{US-CT}^M \times W_V^M \quad (12)$$

The W_Q^M , W_K^M , and W_V^M are trainable parameters. 616
Through the multiplication of Q^M and K^M , an atten- 617
tion matrix emerges, encapsulating the priority of the 618
two modalities within separate nodes. 619

$$A^M = \text{softmax} \left(\frac{Q^M \times K^{MT}}{\sqrt{d_k}} \right) \quad (13)$$

Subsequently, this attention matrix is employed to 620
adjust the relative significance of the two modalities: 621

$$f^M = V^M + A^M \times V^M \quad (14)$$

In the nodal correlation branch, US and CT features 622
are first merged along the dimensional axis: 623

$$f_{US-CT}^N = \text{concat} \left(f_{US}, f_{CT} \right) \quad (15)$$

Then Q^N , K^N , and V^N are obtained respectively: 624

$$Q^N = f_{US-CT}^N \times W_Q^N \quad (16)$$

$$K^N = f_{US-CT}^N \times W_K^N \quad (17)$$

$$V^N = f_{US-CT}^N \times W_V^N \quad (18)$$

The W_Q^N , W_K^N , and W_V^N are trainable parameters. 625
The attention matrix is obtained and employed to 626
delineate the interrelation between distinct nodules: 627

$$A^N = \text{softmax} \left(\frac{Q^N \times K^{NT}}{\sqrt{d_k}} \right) \quad (19)$$

Refined features considering the similarity of dif- 628
ferent nodules emerge: 629

$$f^N = V^N + A^N \times V^N \quad (20)$$

The features from the three branches undergo 630
element-wise multiplication, resulting in the ultimate 631
fused features: 632

$$f_F = f^D \odot f^M \odot f^N \quad (21)$$

633 Nodule boundary extraction in US images

634 Firstly, the nodule boundary width (d) was deter- 635
mined as a multiple (f) of the square region of inter- 636
est's length. For our study, f was set to 0.08. Sec- 637
ondly, the annotated thyroid nodule mask underwent 638
dilation and erosion operations to yield R_{dilation} and 639
 R_{erosion} , respectively, with a kernel size of 3×3 and 640
iterations of $0.5d$. Finally, the nodule's boundary 641
was obtained as the difference between R_{dilation} and 642
 R_{erosion} ($R_{\text{dilation}} - R_{\text{erosion}}$).

643 **Boundaries of surrounding tissue** 644 **extraction in CT images**

645 Firstly, bilateral filtering[36] was applied to preserve
646 the edges while reducing noise. The diameter of the
647 pixel field was set to 7, and the sigma values for both
648 the color space and coordinate space were set to 100.
649 Secondly, the Canny algorithm[37] was employed to
650 further extract the boundaries of the surrounding
651 tissues. The lower and upper threshold values were
652 set to -100 and 200, respectively.

653 **Training Configuration**

654 The base learning rate in our study was set to
655 1×10^{-4} , and we employed a cosine learning rate
656 schedule during the training process. The batch size
657 was set to 30, and we utilized the Adam optimizer
658 to optimize our model. A weight decay of 1×10^{-5}
659 was applied to mitigate overfitting. In this study, a
660 multi-task strategy was employed to address differ-
661 ent tasks. For the classification task, specifically the
662 prediction of the LNM status, we utilized a binary
663 cross-entropy loss function. As for the segmenta-
664 tion tasks, a combination of binary cross-entropy loss
665 and Intersection over Union (IOU) loss functions was
666 utilized. The model was initially trained for the seg-
667 mentation tasks for the first 100 epochs, and then
668 the classification task was added and trained for the
669 remaining 200 epochs.

670 **Interpretability Analysis Methods**

671 We employed the integrated gradients[35] method to
672 enhance the interpretability of our model. Integrated
673 gradients is a feature attribution technique that cal-
674 culates the integral of gradients along the path from
675 a chosen baseline to the input, resulting in an attri-
676 bution value for each input feature. In our study, the
677 baseline is manually specified, and we select a base-
678 line where the predicted probability of our trained
679 model is close to 0.5, indicating equal probabilities
680 for both LNM presence and absence. To determine
681 the contributions of US and CT images, we sum the
682 attributions of each pixel in the respective images. By
683 visualizing the attribution of each pixel, we generate
684 saliency maps for US and CT images.

685 **Statistical Analysis**

686 We assessed the performance of our model using
687 several evaluation metrics, including accuracy, area
688 under the curve (AUC), specificity, sensitivity (also
689 known as recall), precision, and F1-score. To ana-
690 lyze the model's performance across different thresh-
691 olds, we constructed receiver operating characteristic
692 (ROC) curves, plotting sensitivity against specificity.

693 **Hardware and Software**

694 The computational resources utilized include an Intel
695 10900K CPU with a clock speed of 3.7GHz and 20
696 threads. The graphics card employed is a GEFORCE
697 RTX 3090, equipped with 10752 CUDA cores and
698 24GB of graphics memory. The programming lan-
699 guage used for implementation is Python 3.9.7, and
700 the deep learning framework employed is PyTorch
701 1.10.0.

702 **Data availability**

703 Though this study was carried out with participant
704 consent, the dataset remains restricted in public ac-
705 cess. For research inquiries, de-identified data can
706 be obtained from the corresponding author upon
707 reasonable request.

708 **Code availability**

709 The code for model development and
710 interpretability analysis is accessible at
711 <https://github.com/li10107/DGFNet>.

712 **References**

- 713 [1] YuJiao Deng et al. "Global burden of thyroid
714 cancer from 1990 to 2017". In: *JAMA network
715 open* 3.6 (2020), e208759–e208759.
- 716 [2] Hyuna Sung et al. "Global cancer statistics
717 2020: GLOBOCAN estimates of incidence and
718 mortality worldwide for 36 cancers in 185 coun-
719 tries". In: *CA: a cancer journal for clinicians* 71.3
720 (2021), pp. 209–249.
- 721 [3] Kyu Eun Lee et al. "Ipsilateral and contralateral
722 central lymph node metastasis in papillary thy-
723 roid cancer: patterns and predictive factors of
724 nodal metastasis". In: *Head & neck* 35.5 (2013),
725 pp. 672–676.
- 726 [4] David T Hughes and Gerard M Doherty. "Cen-
727 tral neck dissection for papillary thyroid can-
728 cer". In: *Cancer Control* 18.2 (2011), pp. 83–88.
- 729 [5] Bryan R Haugen et al. "2015 American Thy-
730 roid Association management guidelines for
731 adult patients with thyroid nodules and differ-
732 entiated thyroid cancer: the American Thyroid
733 Association guidelines task force on thyroid
734 nodules and differentiated thyroid cancer". In:
735 *Thyroid* 26.1 (2016), pp. 1–133.

- 736 [6] Claudio Gambardella et al. "The role of prophylactic central compartment lymph node dissection in elderly patients with differentiated thyroid cancer: a multicentric study". In: *BMC surgery* 18.1 (2019), pp. 1–8. 787
- 737 788
- 738 789
- 739 790
- 740
- 741 [7] Mostafa Alabousi et al. "Diagnostic test accuracy of ultrasonography vs computed tomography for papillary thyroid cancer cervical lymph node metastasis: A systematic review and meta-analysis". In: *JAMA Otolaryngology–Head & Neck Surgery* 148.2 (2022), pp. 107–118. 791
- 742 792
- 743 793
- 744 794
- 745 795
- 746 796
- 747 [8] Yinlong Yang et al. "Prediction of central compartment lymph node metastasis in papillary thyroid microcarcinoma". In: *Clinical endocrinology* 81.2 (2014), pp. 282–288. 797
- 748 798
- 749 799
- 750
- 751 [9] Jiang Zhu et al. "Application of machine learning algorithms to predict central lymph node metastasis in T1-T2, non-invasive, and clinically node negative papillary thyroid carcinoma". In: *Frontiers in medicine* 8 (2021), p. 635771. 800
- 752 801
- 753 802
- 754 803
- 755 804
- 756 [10] Jong-Lyel Roh, Jin-Man Kim, and Chan Il Park. "Central lymph node metastasis of unilateral papillary thyroid carcinoma: patterns and factors predictive of nodal metastasis, morbidity, and recurrence". In: *Annals of surgical oncology* 18 (2011), pp. 2245–2250. 805
- 757 806
- 758 807
- 759 808
- 760 809
- 761 810
- 762 [11] Meng Jiang et al. "Nomogram based on shear-wave elastography radiomics can improve pre-operative cervical lymph node staging for papillary thyroid carcinoma". In: *Thyroid* 30.6 (2020), pp. 885–897. 811
- 763 812
- 764 813
- 765 814
- 766 815
- 767 [12] Jinhua Yu et al. "Lymph node metastasis prediction of papillary thyroid carcinoma based on transfer learning radiomics". In: *Nature communications* 11.1 (2020), p. 4807. 816
- 768 817
- 769 818
- 770 819
- 771 [13] Tongtong Liu et al. "Comparison of the application of B-mode and strain elastography ultrasound in the estimation of lymph node metastasis of papillary thyroid carcinoma based on a radiomics approach". In: *International journal of computer assisted radiology and surgery* 13 (2018), pp. 1617–1627. 820
- 772 821
- 773 822
- 774 823
- 775 824
- 776 825
- 777 826
- 778 [14] Vivian Y Park et al. "Radiomics signature for prediction of lateral lymph node metastasis in conventional papillary thyroid carcinoma". In: *PLoS One* 15.1 (2020), e0227315. 827
- 779 828
- 780 829
- 781 830
- 782 831
- 783 832
- 784 [15] Jingjing Li et al. "Computed tomography-based radiomics model to predict central cervical lymph node metastases in papillary thyroid carcinoma: a multicenter study". In: *Frontiers in Endocrinology* 12 (2021), p. 741698. 833
- 785 834
- 786 835
- 787 836
- 788 837
- [16] Yun Peng et al. "Prediction of central lymph node metastasis in cN0 papillary thyroid carcinoma by CT radiomics". In: *Academic Radiology* 30.7 (2023), pp. 1400–1407. 838
- [17] Shanshan Zhao et al. "Combined Conventional Ultrasound and Contrast-Enhanced Computed Tomography for Cervical Lymph Node Metastasis Prediction in Papillary Thyroid Carcinoma". In: *Journal of Ultrasound in Medicine* 42.2 (2023), pp. 385–398. 839
- [18] Jana Lipkova et al. "Artificial intelligence for multimodal data integration in oncology". In: *Cancer cell* 40.10 (2022), pp. 1095–1110. 840
- [19] Kevin M Boehm et al. "Harnessing multimodal data integration to advance precision oncology". In: *Nature Reviews Cancer* 22.2 (2022), pp. 114–126. 841
- [20] Julián N Acosta et al. "Multimodal biomedical AI". In: *Nature Medicine* 28.9 (2022), pp. 1773–1784. 842
- [21] Kevin M Boehm et al. "Multimodal data integration using machine learning improves risk stratification of high-grade serous ovarian cancer". In: *Nature cancer* 3.6 (2022), pp. 723–733. 843
- [22] Xuejun Qian et al. "Prospective assessment of breast cancer risk from multimodal multi-view ultrasound images via clinically applicable deep learning". In: *Nature biomedical engineering* 5.6 (2021), pp. 522–532. 844
- [23] Richard J Chen et al. "Pan-cancer integrative histology-genomic analysis via multimodal deep learning". In: *Cancer Cell* 40.8 (2022), pp. 865–878. 845
- [24] Hong-Yu Zhou et al. "A transformer-based representation-learning model with unified processing of multimodal input for clinical diagnostics". In: *Nature Biomedical Engineering* (2023), pp. 1–13. 846
- [25] Ye Zhu et al. "Vision+ X: A Survey on Multimodal Learning in the Light of Data". In: *arXiv preprint arXiv:2210.02884* (2022). 847
- [26] Chao Zhang et al. "Multimodal intelligence: Representation learning, information fusion, and applications". In: *IEEE Journal of Selected Topics in Signal Processing* 14.3 (2020), pp. 478–493. 848
- [27] Hassan Akbari et al. "Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text". In: *Advances in Neural Information Processing Systems* 34 (2021), pp. 24206–24221. 849

- 838 [28] Kaiming He et al. “Masked autoencoders are
839 scalable vision learners”. In: *Proceedings of the*
840 *IEEE/CVF conference on computer vision and pat-*
841 *tern recognition*. 2022, pp. 16000–16009.
- 842 [29] Xiaohan Wang, Linchao Zhu, and Yi Yang.
843 “T2vlad: global-local sequence alignment for
844 text-video retrieval”. In: *Proceedings of the*
845 *IEEE/CVF Conference on Computer Vision and*
846 *Pattern Recognition*. 2021, pp. 5079–5088.
- 847 [30] Diping Song et al. “Deep relation transformer
848 for diagnosing glaucoma with optical coher-
849 ence tomography and visual field function”. In:
850 *IEEE Transactions on Medical Imaging* 40.9 (2021),
851 pp. 2392–2402.
- 852 [31] Shuai Zheng et al. “Multi-modal graph learn-
853 ing for disease prediction”. In: *IEEE Transac-*
854 *tions on Medical Imaging* 41.9 (2022), pp. 2207–
855 2216.
- 856 [32] Richard J Chen et al. “Multimodal co-attention
857 transformer for survival prediction in gigapixel
858 whole slide images”. In: *Proceedings of the*
859 *IEEE/CVF International Conference on Computer*
860 *Vision*. 2021, pp. 4015–4025.
- 861 [33] Tsai Hor Chan et al. “Histopathology Whole
862 Slide Image Analysis With Heterogeneous
863 Graph Representation Learning”. In: *Proceed-*
864 *ings of the IEEE/CVF Conference on Computer*
865 *Vision and Pattern Recognition*. 2023, pp. 15661–
866 15670.
- 867 [34] Kaiming He et al. “Deep residual learning for
868 image recognition”. In: *Proceedings of the IEEE*
869 *conference on computer vision and pattern recogni-*
870 *tion*. 2016, pp. 770–778.
- 871 [35] Mukund Sundararajan, Ankur Taly, and Qiqi
872 Yan. “Axiomatic attribution for deep networks”.
873 In: *International conference on machine learning*.
874 PMLR. 2017, pp. 3319–3328.
- 875 [36] Carlo Tomasi and Roberto Manduchi. “Bilateral
876 filtering for gray and color images”. In: *Sixth*
877 *international conference on computer vision (IEEE*
878 *Cat. No. 98CH36271)*. IEEE. 1998, pp. 839–846.
- 879 [37] John Canny. “A computational approach to
880 edge detection”. In: *IEEE Transactions on pat-*
881 *tern analysis and machine intelligence* 6 (1986),
882 pp. 679–698.

883 Acknowledgments

884 This work was supported by the National Natural
885 Science Foundation of China (No. 82071946), the
886 Natural Science Foundation of Zhejiang Province
887 (No. LZ Y21F030001), the Pioneer and Leading Goose

R&D Program of Zhejiang (No. 2023C04039), the
888 National Key Research and Development Program
889 of China (2022YFF0608403), Youth Research Fund
890 Project of Shaoxing People’s Hospital (Grant Num-
891 ber 2022YB07), and the fund of Zhejiang Province
892 Medical and Health Science and Technology Project
893 (No. 2023KY581). We thank Y.G. for providing us
894 external validation set.
895

896 Author contributions

M.H., C.S., X.L., and D.X. conceived and planned
897 the study. C.S., G.L., and X.L. designed the research
898 framework. J.Y., C.P., S.Z., and J.Y. collected the
899 raw US and CT images, patients’ clinical informa-
900 tion, and image annotation. G.L., Y.H., and X.F. per-
901 formed the data preprocessing and conducted the
902 performance analysis. G.L. designed the multimodal
903 fusion method and carried out model interpretation
904 analysis. G.L. and C.S. wrote the manuscript. All
905 authors commented on the manuscript.
906

907 Competing interests

The authors declare no competing interests.
908

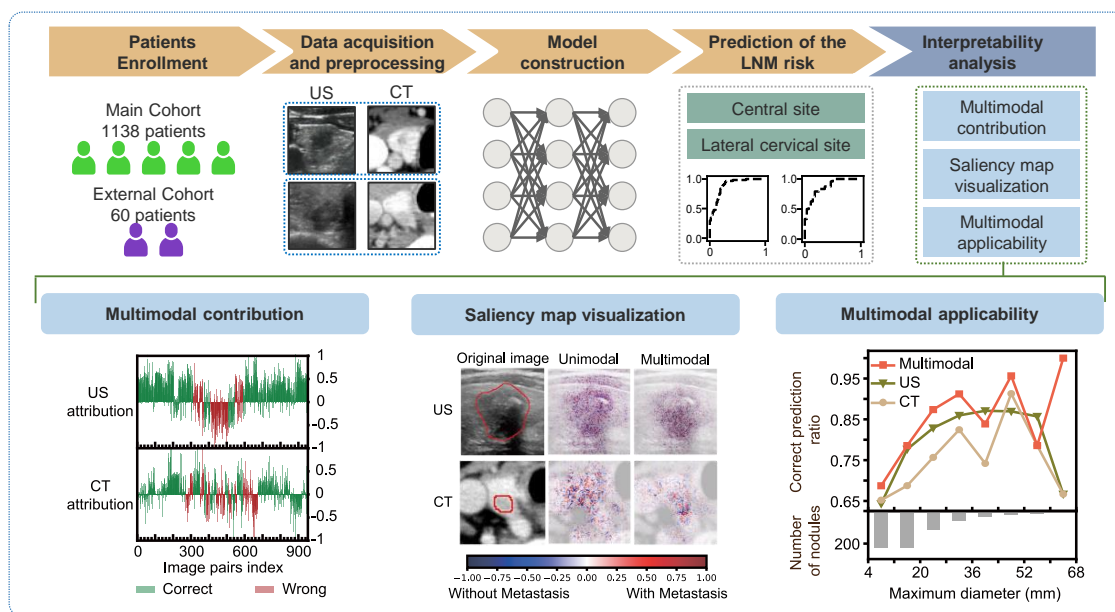


Figure 1: Overall AI system for LNM risk prediction. The main cohort was employed for AI system development and evaluation, while the external cohort assessed the system's generalizability. After preprocessing, paired US and CT images are input into DGFNet, our deep learning model, to predict LNM status in central and lateral cervical regions. Post-AI system development, we conducted an extensive interpretability analysis comprising multimodal contribution assessment, saliency map visualization, and multimodal applicability evaluation.

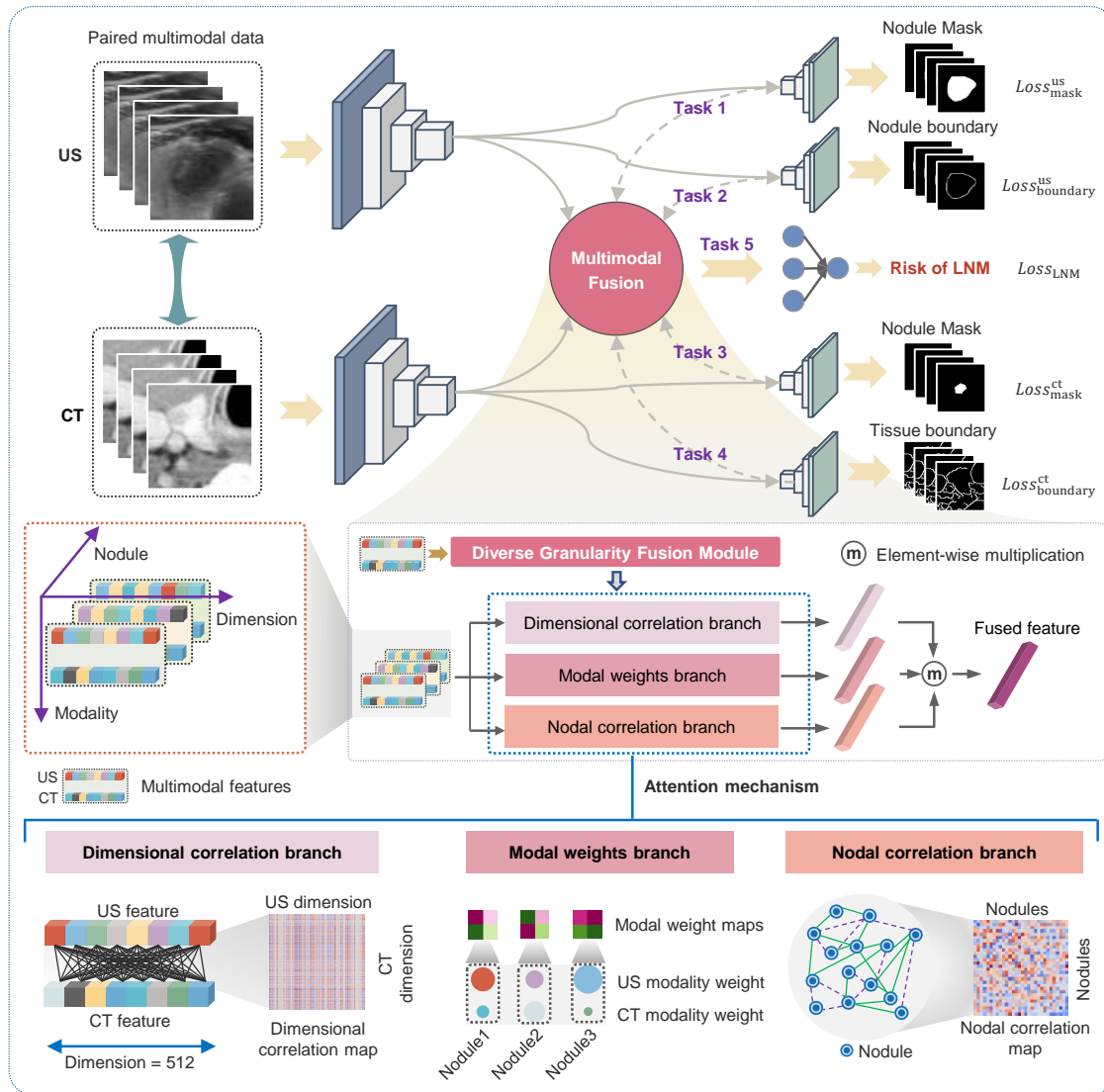


Figure 2: DGFNet architecture. DGFNet consists of three branches: the US branch, CT branch, and multimodal branch. Each US and CT branch incorporates an encoder and two decoders. DGFNet concurrently performs five tasks: nodal mask and boundary segmentation in US images (guiding the model to focus on internal and marginal nodule features), boundary segmentation of nodules and surrounding tissues in CT images (guiding the model to focus on nodule and surrounding tissue features in CT images), and the final LNM prediction. The fusion of multimodal features in the latent space occurs within the diverse granularity fusion module, and the final results are generated by subsequent fully connected layers. The diverse granularity fusion module includes the dimensional correlation branch, modal weights branch, and nodal correlation branch, amalgamating characteristics from both modalities to provide a diverse granularity information integration. A detailed explanation of this module is available in the Methods section.

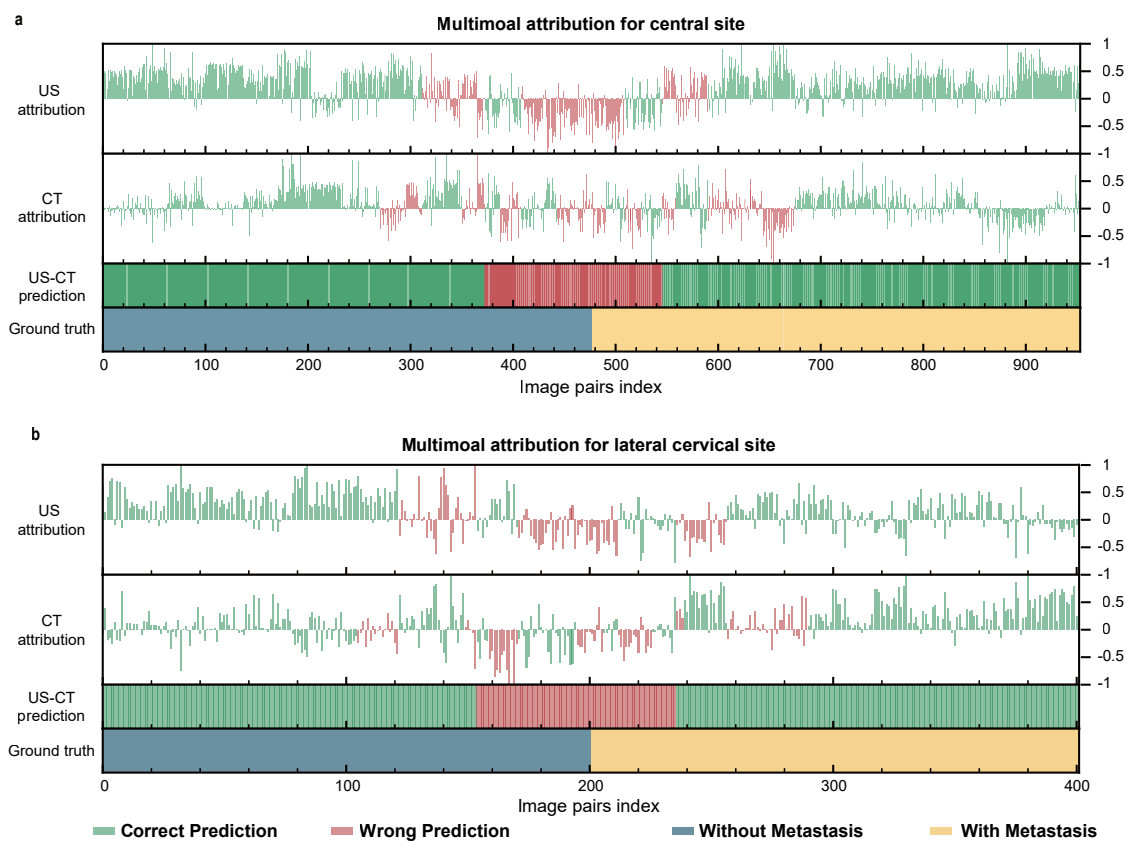


Figure 3: Attribution analysis of US and CT in predicting LNM status at central (a) and lateral cervical (b) sites. Each subfigure comprises four panels, with the shared horizontal axes indicating nodule indices. The central site includes 954 nodules, while the lateral cervical site includes 402 nodules. The values of the top two panels display attributions from US and CT images in the multimodal prediction, respectively. Column colors denote unimodal predictions, where green signifies accurate predictions and red indicates inaccuracies. The third panel illustrates the multimodal prediction results. Panel 4 represents the ground truth.

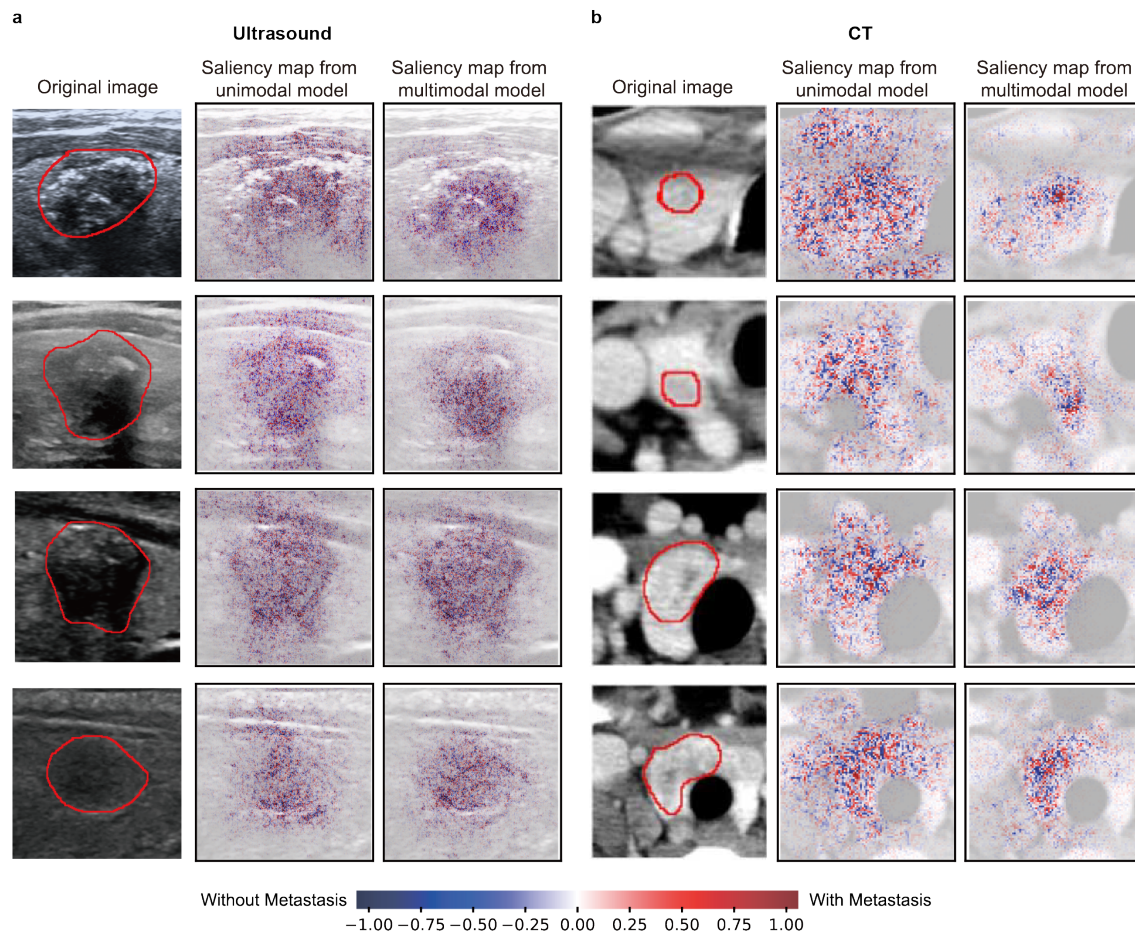


Figure 4: Examples of saliency map visualization results at central site. In these instances, both the US and CT unimodal models initially generated inaccurate predictions, whereas the multimodal models effectively rectified these to provide accurate predictions. The red curve delineates the nodule's boundary in the original US and CT images. The color red signifies an elevated likelihood of LNM development, whereas the color blue signifies the contrary.

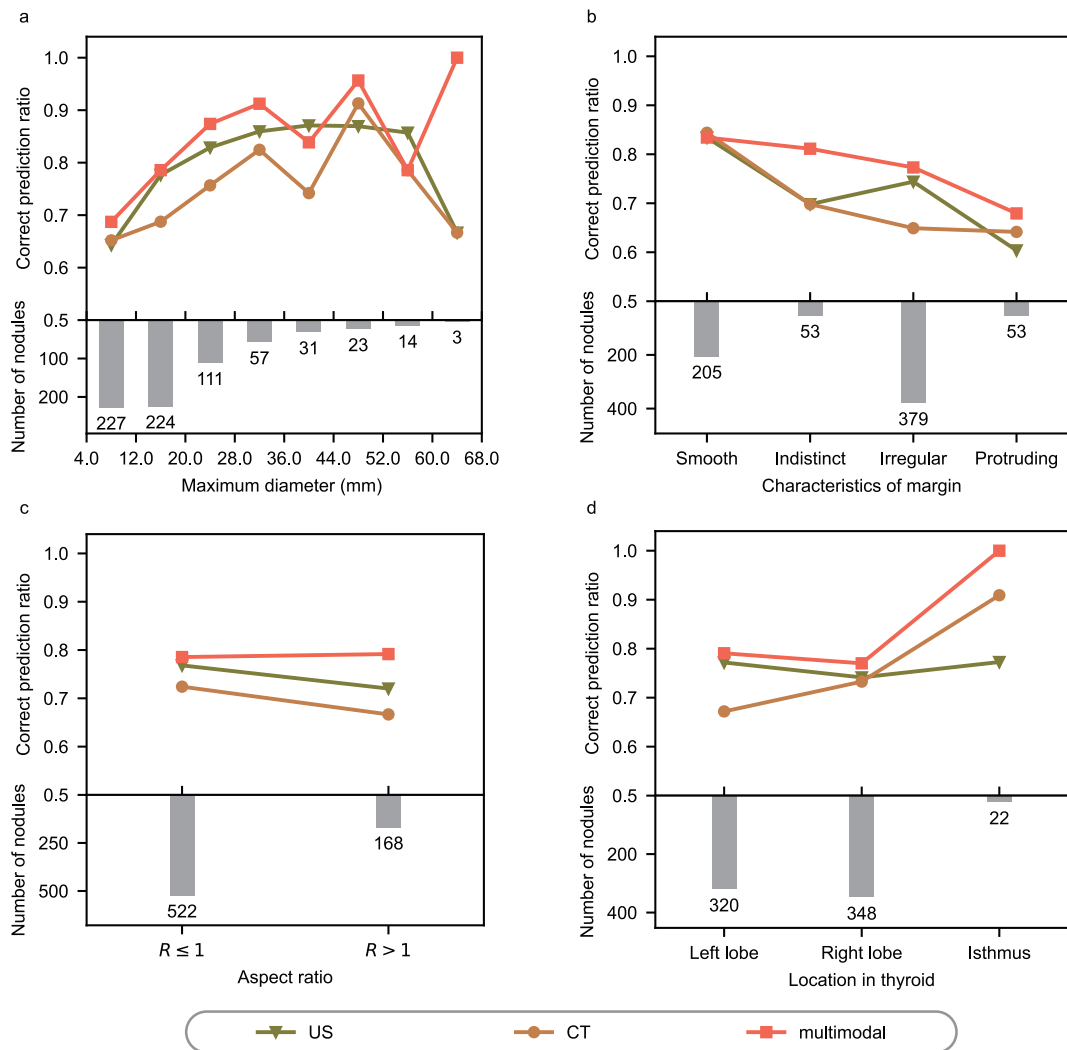


Figure 5: Distribution of nodules with varied attributes and associated correct predictions ratio in central Site. Attributes encompass nodal maximum diameter (considering the larger of the maximum diameters from transverse and longitudinal US views) in US image(a), characteristics of margin (b), aspect ratio (calculated as the height divided by the width in transverse views) of nodules(c), and location in thyroid (d).