

# *BiomedParse*: a biomedical foundation model for image parsing of everything everywhere all at once

Theodore Zhao<sup>1\*</sup>¶, Yu Gu<sup>1\*</sup>, Jianwei Yang<sup>1</sup>, Naoto Usuyama<sup>1</sup>, Ho Hin Lee<sup>1</sup>, Tristan Naumann<sup>1</sup>, Jianfeng Gao<sup>1</sup>, Angela Crabtree<sup>2</sup>, Jacob Abel<sup>2</sup>, Christine Moung-Wen<sup>2</sup>, Brian Piening<sup>2</sup>, Carlo Bifulco<sup>2</sup>, Mu Wei<sup>1,‡</sup>, Hoifung Poon<sup>1,‡§</sup>, Sheng Wang<sup>3,‡</sup>

<sup>1</sup>Microsoft Research, Redmond, WA, USA

<sup>2</sup>Providence Genomics, Portland, OR, USA

<sup>3</sup>Paul G. Allen School of Computer Science and Engineering, University of Washington, Seattle, WA, USA

<https://aka.ms/biomedparse-project>

## Abstract

Biomedical image analysis is fundamental for biomedical discovery in cell biology, pathology, radiology, and many other biomedical domains. Holistic image analysis comprises interdependent subtasks such as segmentation, detection, and recognition of relevant objects. Traditionally, these tasks are tackled separately. For example, there have been a lot of works focusing on segmentation alone, completely ignoring key semantic information in downstream tasks of detection and recognition. In contrast, image parsing is a unifying framework that jointly pursues these tasks by leveraging their interdependencies such as the semantic label of a segmented object. Here, we propose *BiomedParse*, a biomedical foundation model for imaging parsing that can jointly conduct segmentation, detection, and recognition for 82 object types across 9 imaging modalities. Through joint learning, we can improve accuracy for individual tasks and enable novel applications such as segmenting all relevant objects in an image through a text prompt, rather than requiring users to laboriously specify the bounding box for each object. Interestingly, we can train *BiomedParse* using no more than standard segmentation datasets. The key is to leverage readily available natural-language labels or descriptions accompanying those datasets and use GPT-4 to harmonize the noisy, unstructured text information with established biomedical object ontologies. We created a large dataset comprising over six million triples of image, segmentation mask, and textual description. On image segmentation, we showed that *BiomedParse* is broadly applicable, outperforming state-of-the-art methods on 102,855 test image-mask-label triples across 9 imaging modalities (*everything*). *BiomedParse* is also able to identify invalid user inputs describing objects that do not exist in the image. On object detection, which aims to locate a specific object of interest, *BiomedParse* again attained state-of-the-art performance, especially on objects with irregular shapes (*everywhere*). On object recognition, which aims to identify all objects in a given image along with their semantic types, we showed that *BiomedParse* can simultaneously segment and label all biomedical objects in an image (*all at once*). In summary, *BiomedParse* is an all-in-one tool for biomedical image analysis by jointly solving segmentation, detection, and recognition. It is broadly applicable to all major biomedical image modalities, paving the path for efficient and accurate image-based biomedical discovery.

---

\* Joint first authors

† Project lead

¶ Main technical contribution

‡ Corresponding authors: muhsin.wei@microsoft.com, swang@cs.washington.edu, hoifung@microsoft.com

§ Lead contact

## Introduction

Biomedical image analysis is critical to biomedical discovery because imaging is one of the most important tools for studying physiology, anatomy, and function at multiple scales from the organelle level to the organ level [1, 2]. Holistic image analysis comprises multiple subtasks, such as segmentation, detection, and recognition of biomedical objects. Segmentation aims to divide an image into segments representing different objects, often requiring the aid of a user-provided bounding box for each object of interest [3, 4]. Detection aims to identify the location of an object of interest in the image [5], whereas recognition aims to identify all objects within an image [6]. Standard image analysis methods typically approach these tasks separately, using specialized tools for individual tasks [7]. Despite their encouraging performance, such a disjoint approach misses significant opportunities for joint learning and reasoning across these interdependent tasks.

For example, a lot of prior image analysis works focus on segmentation alone, thus ignoring key semantic information from downstream tasks of object detection and recognition. This results in sub-optimal segmentation while imposing substantial burden on users, as many state-of-the-art segmentation tools require users to provide a tight bounding box indicating the location of an object of interest [8, 9]. The bounding-box requirement leads to three limitations. First, users have to manually draw bounding boxes in the image, which requires domain expertise to identify the locations and shapes of the target objects. Second, bounding boxes, which are often rectangular, fall short of accurately representing objects with irregular or complex shapes. Third, bounding box-based approaches are not scalable for images containing a large number of objects, such as segmenting cells in a whole-slide pathology image, since users need to provide a bounding box for each object.

In this paper, we propose to approach biomedical image analysis as *image parsing*, a unifying framework for joint learning and reasoning across segmentation, detection, and recognition [10–12]. Specifically, we have developed *BiomedParse*, a biomedical foundation model for image parsing that is capable of carrying out all three tasks by leveraging their interdependencies, thus addressing key limitations in traditional methods. In particular, joint learning of object detection and recognition eliminates the need for user-specified bounding boxes, as segmentation can be done using semantic labels from text prompt alone.

The major bottleneck for pretraining *BiomedParse* is data. While biomedical segmentation datasets abound [13–15], there are relatively few prior works on object detection and recognition in biomedicine, let alone datasets covering all three tasks. To address this problem, we propose a novel approach for pretraining *BiomedParse* using no more than standard segmentation datasets. The key insight is to leverage readily available natural-language labels or descriptions accompanying those datasets and use GPT-4 to harmonize these noisy, unstructured texts with established biomedical object ontologies. This enables us to construct *BiomedParseData*, a biomedical image parsing dataset comprising 3.4 million triples of image, segmentation mask, and semantic label of the biomedical object and 6.8 million image-mask-description triples, from over 1 million images. The semantic labels encompass 82 major biomedical object types across 9 imaging modalities.

Unlike segmentation methods that focus on identifying salient segment boundary within a bounding box, *BiomedParse* learns to model typical shape of each object class, thus mimicking how humans perceive objects in an image. *BiomedParse* can segment images using text prompts alone (e.g. “inflammatory cells in breast pathology”), without requiring any user-specified localization such as bounding boxes. Consequently, *BiomedParse* can better recognize and segment objects of irregular and complex shapes, which are very challenging for traditional methods using rectangular bounding boxes. Moreover, *BiomedParse* can recognize all objects in an image, without requiring any user input prompt.

We conduct a large-scale study to evaluate *BiomedParse* on 102,855 held-out image-mask-label triples across 9 modalities for segmentation, detection, and recognition. On segmentation, *BiomedParse* established new state-of-the-art results, outperforming prior best methods such as MedSAM [9] and SAM [8].

Moreover, using text prompts alone, *BiomedParse* is much more scalable than these prior methods that require orders of magnitude more user operations in specifying object-specific bounding boxes to perform competitively. We also demonstrated that *BiomedParse* can accurately detect invalid text prompts describing non-existent objects in the image. On detection, we show that *BiomedParse* learns accurate modeling of all object classes, including those with irregular shapes. This results in even larger improvement in image analysis accuracy for such objects, attaining a 0.857 Dice score that is 39.6% higher than the best-competing method. On recognition, we show how *BiomedParse* can accurately segment and label all objects without any user-specified text prompt. Collectively, we introduce a biomedical foundation model for image parsing, achieving superior performance on segmentation, detection, and recognition, paving the path for large-scale image-based biomedical discovery.

## Results

### Overview of *BiomedParse* and *BiomedParseData*

To develop a model that can jointly conduct segmentation, detection, and recognition, we need a supervision dataset that covers all three tasks. To our best knowledge, no such a dataset exists. To this end, we created the first such dataset *BiomedParseData* by combining 45 biomedical image segmentation datasets and using GPT-4 to generate the canonical semantic label for each segmented object.

The key insight is that existing segmentation datasets often contain valuable semantic information about the segmented objects. However, such information typically resides in noisy and inconsistent natural-language text descriptions that do not conform to standard biomedical ontologies. To address this challenge, we use GPT-4 to create a unifying biomedical object ontology for image analysis and harmonize natural-language descriptions with this ontology (see **Methods**). This ontology encompasses three main categories (histology, organ, abnormality), 15 meta-object types, and 82 specific object types (**Fig. 1a**).

The resulting *BiomedParseData* contains 3.4 million distinct image-mask-label triples, spanning 9 imaging modalities and 25 anatomic sites (**Fig. 1b, Supplementary Figure 2**), representing a large-scale and diverse dataset for semantic-based biomedical image analysis.

To make *BiomedParse* better equipped in handling diverse text prompts not covered by the canonical semantic labels, we also use GPT-4 to synthesize synonymous text descriptions for each semantic label and sample from them during training. This yielded a total of 6.8 million image-mask-description triples (see **Methods, Supplementary Figure 1** and 3).

While our method does not use bounding boxes, prior state-of-the-art methods such as MedSAM and SAM generally require pre-specified bounding boxes. We consider two scenarios to provide the bounding boxes: oracle bounding box (the minimum rectangular bounding box covering a segmented object) and bounding box created by Grounding DINO [16], a state-of-the-art object detection method that can generate bounding boxes from text prompt of an object label. (Grounding DINO does not perform segmentation.)

*BiomedParse* adopts a modular design under the SEEM architecture [56], comprising an image encoder (for encoding the input image), a text encoder (for encoding the text prompt), a mask decoder (for outputting segmentation mask), and a meta-object classifier (for joint training of image encoder with object semantics). See **Fig. 1c**. The image and text encoders were initialized using state-of-the-art Focal [17] and PubMedBERT [18], respectively.

Before evaluating image analysis results, we first examine the quality of embeddings derived from *BiomedParse*. Specifically, we compare the text embeddings from *BiomedParse* with those from PubMedBERT. We found that embeddings from *BiomedParse* can better distinguish fine-grained cell types, with a Silhouette score of 0.89 that is much higher than using the embeddings from PubMedBERT (**Fig. 1d, Supplementary Figure 4**). We also compare the image embeddings from *BiomedParse* with those from Focal. We observed that embeddings from *BiomedParse* are more predictive of tumor malignancy on a pathology

dataset [19] (**Fig. 1e**). The superior performance of the text and image embeddings from *BiomedParse* necessitates the training of *BiomedParse* using *BiomedParseData*, raising our confidence that *BiomedParse* can be an effective approach for biomedical image analysis.

### Accurate and scalable segmentation across nine modalities

We first evaluated *BiomedParse* on biomedical image segmentation using the held-out set comprising 102,855 test instances (image-mask-label triples) across 9 imaging modalities (**Fig. 2a, Supplementary Figure 4**). We observed that *BiomedParse* achieved the best Dice score, even against the best-competing method MedSAM with oracle bounding box as input (paired t-test  $p$ -value  $< 10^{-4}$ ). In the more realistic setting when MedSAM or SAM is supplied with bounding boxes generated by Grounding DINO, the superiority of *BiomedParse* is even more prominent in end-to-end biomedical object detection and segmentation, especially in more challenging modalities such as pathology and CT where irregular-shaped objects abound. By training on domain-specific datasets, both *BiomedParse* and MedSAM outperform general-domain methods such as SAM. We showed examples comparing *BiomedParse* segmentation and the ground truth across multiple imaging modalities, demonstrating the generalizability of *BiomedParse* (**Fig. 2b**). We further compared *BiomedParse* on a benchmark created by MedSAM [9] encompassing 50 tasks and again observed the best performance by *BiomedParse*, even against MedSAM with oracle bounding box (paired t-test  $p$ -value  $< 10^{-2}$ ), further demonstrating the superiority of *BiomedParse* (**Supplementary Figure 5**).

In addition to being more accurate, *BiomedParse* is more scalable compared to bounding box-based approaches, which stems from the generalizability of text prompts across images of the same modality or anatomical site, thus eliminating the need for laborious user operations in providing a tight bounding box for each object. To demonstrate this, we compared *BiomedParse* and prior state-of-the-art methods MedSAM and SAM on a cell segmentation dataset with 42 colon pathology images (**Fig. 2c**). Using a single text prompt “glandular structure in colon pathology image”, *BiomedParse* achieves a 0.942 median Dice score, whereas neither SAM nor MedSAM achieves a median Dice score higher than 0.75 without tight bounding boxes as input. In fact, to achieve competitive results comparable to *BiomedParse* with a single text prompt, MedSAM requires the users to supply a tight bounding box for each of the 430 cells in these images (**Fig. 2c**). In general, our results reveal that the bounding box-based approach is much less accurate on irregular-shaped objects, such as tumors and abnormal cells (**Fig. 2d,e**). In contrast, *BiomedParse* still attained highly accurate segmentation for such objects. The scalability and accuracy of *BiomedParse* bode well for its utility in real-world applications.

*BiomedParse* can also detect invalid text prompts (e.g., the request to identify a brain tissue in a chest X-Ray image), by calculating a  $p$ -value using Kolmogorov–Smirnov (K-S) test (see **Methods**). From preliminary experiments, we found that invalid text prompts have an average K-S test  $p$ -value smaller than  $10^{-3}$  while the valid ones have an average K-S test  $p$ -value above 0.1 (**Fig. 2f**). Using 0.01 as the  $p$ -value cutoff, *BiomedParse* can achieve an estimated performance of 0.93 precision and 1.00 recall on detecting invalid input (**Fig. 2g**). *BiomedParse* substantially outperformed Grounding DINO on invalid input detection (AUROC 0.99 vs 0.61). See **Fig. 2h,i**. This enables *BiomedParse* to perform recognition by enumerating candidate object types in the ontology, skipping invalid text prompts and generating segmentation masks for valid object labels.

### Accurate detection of irregular-shaped objects

Next, we evaluated the performance of *BiomedParse* on object detection, where the model is asked to identify an object of interest in the image. *BiomedParse* can resolve natural-language variations and accept text prompts that do not exactly match any semantic label in the ontology. In the previous section, we already show that *BiomedParse* outperformed bounding-box-based methods in general. Additionally, since *Biomed-*

*Parse* learns semantic representation for individual object types. We hypothesize that its superiority over prior methods will be even more pronounced in detecting irregular-shaped objects. To verify this, we show the aggregate attention map of each object type learned by *BiomedParse* on test images unseen during training and observed that they faithfully reflect object shapes, including many irregular-shaped objects (**Fig. 3a**). Next, we define three metrics to assess the regularity of an object, including *convex ratio* (i.e., the ratio of the object size to the tightest convex size), *box ratio* (i.e., the ratio of the object size to the tightest rectangle size), and *rotational inertia* (i.e., the difficulty in changing the rotational velocity) (see **Methods**). We found that the improvements of *BiomedParse* over SAM and MedSAM are strongly correlated with these metrics (average correlation 0.829), indicating that our method has a larger improvement on irregular-shaped objects (**Fig. 3b-d, Supplementary Figure 6**). **Fig. 3e** illustrates a few examples comparing *BiomedParse* and MedSAM on detecting irregular-shaped objects. Furthermore, we show that *BiomedParseData* has higher average object irregularity than the datasets used by MedSAM (**Fig. 3f,g, Supplementary Figure 7**), and the improvement of *BiomedParse* is also larger on *BiomedParseData* (**Fig. 3h**), highlighting the benefit from joint learning of object semantics in detecting the more challenging irregular-shaped objects.

## Object recognition using the segmentation ontology

In our final analysis, we explore *BiomedParse*'s capacity for object recognition, which aims to simultaneously segment and label every object within an image. Provided with an image, along with its modality and anatomical site, *BiomedParse* iteratively performs detection and segmentation for all candidate object types within the ontology of that modality and anatomical site, and the segmented masks are aggregated to ensure spatial cohesion among adjacent pixels (see **Methods**). This enables *BiomedParse* to accurately conduct object recognition, as evidenced in **Fig. 4a**, where objects are accurately identified and segmented with an average Dice score of 0.94.

Grounding DINO [16] is the state-of-the-art general-domain object recognition system but it does not perform segmentation, which makes Grounding DINO and *BiomedParse* not directly comparable. We circumvent this by casting the object recognition task as a binary classification problem: given an input image and a candidate object type, the model determines whether the image contains at least one object of the given type. In this classification formulation, we observed that *BiomedParse* substantially outperformed Grounding DINO with a 25.0%, 87.9%, 74.5% improvement on precision, recall, and F-1, respectively (**Fig. 4b-d**). The improvement over Grounding DINO is even larger when more objects are present in the image (**Fig. 4e**).

Next, we evaluated the performance of *BiomedParse* on end-to-end object recognition using weighted average Dice score. Compared with MedSAM and SAM using Grounding DINO for recognition and bounding box generation, *BiomedParse* outperformed them by a large margin (**Fig. 4f, Supplementary Figure 8**). Similar to our observation on object identification, the improvement over comparison approaches is even larger when more objects are present in the image (**Fig. 4g**). These results indicate *BiomedParse*'s ability to identify all objects in an image, offering an effective tool for holistic image analysis.

Finally, we evaluated *BiomedParse* on real-world data from the Providence Health System (**Fig. 5**). We performed object recognition here by asking *BiomedParse* to identify and segment all relevant cells in the pathology slides. We found that the annotations by *BiomedParse* correctly identified regions of immune cells and cancer cells, attaining high consistency with the pathologist annotations. While pathologists tend to focus on a specific region of cell type and provide coarse-grained annotations, *BiomedParse* can precisely label all relevant cells as specified in the ontology, indicating the potential for *BiomedParse* to help alleviate clinician burdens in real-world clinical applications.

## Discussion

We have presented *BiomedParse*, a biomedical foundation model for image analysis based on image parsing, and a large-scale image parsing dataset *BiomedParseData* containing 3.4 million image-mask-label triples and 6.8 million image-mask-description triples. In contrast to existing biomedical foundational models that require users to provide a tight bounding box for each object to segment, *BiomedParse* is bounding box-free, and can perform holistic image analysis with segmentation, detection, and recognition all at once. We conducted a large-scale evaluation on 102,855 test image-mask-label triples across 9 modalities. *BiomedParse* attained new state-of-the-art results, substantially outperforming prior best methods such as MedSAM and SAM, even when they were equipped with oracle bounding box as input. The improvement is even larger when the objects have irregular shapes or when an image contains a large number of objects. We also validated the accuracy and scalability of *BiomedParse* on previous unseen real-world data from Providence Health System. Collectively, *BiomedParse* offers an accurate, scalable, and robust biomedical image analysis tool that can be broadly applied to various modalities and applications, paving the path for image-based biomedical discovery.

The image analysis field has witnessed rapid development in the past decade. Since its inception in 2015, the U-Net architecture has revolutionized the field of automatic pixel-wise prediction through supervised training [20, 21]. This groundbreaking work laid the foundation for a diverse array of network structures, ranging from advanced convolution-network designs to vision-transformer models [22–37]. Recent advances in image detection and recognition, such as developments in object detection frameworks like Faster R-CNN [38] and YOLOv4 [39], have significantly enhanced capabilities in identifying and localizing anatomical features with high precision. The introduction of SAM marked a significant milestone by demonstrating the model’s ability to generalize segmentation to previously unseen classes, utilizing visual prompts such as points and bounding boxes as guides [8].

Despite the proliferation of advances in the general domain, research on adapting them for large-scale biomedical image analysis across a wide range of organ or tissue classes remains relatively sparse [40]. MedSAM is a notable exception by adapting SAM to the medical realm through continued training on a large number of biomedical segmentation datasets, establishing the state of art in biomedical image analysis. However, like SAM, MedSAM focuses on segmentation alone, thus ignoring valuable semantic information from related tasks of detection and recognition. Consequently, both SAM and MedSAM require users to provide labor-intensive input such as the tight bounding box for each object to segment, which is hard to scale and very challenging for objects with irregular shapes [9].

We propose *BiomedParse* to overcome these challenges. By joint learning across segmentation, detection, and recognition in the unifying framework of image parsing, and by using GPT-4 to harmonize noisy object descriptions, *BiomedParse* was able to acquire novel capabilities such as identifying and segmenting objects of interest using text prompt alone, as well as recognizing all objects in an image by leveraging the segmentation ontology. This represents an important step toward scaling holistic image analysis in biomedicine and real-world clinical applications.

A particularly exciting area for biomedical image analysis is the application in cellular images such as H&E staining and Multiplexed ImmunoFluorescence (MxIF) imaging. This could help elucidate the size, shape, texture, and spatial relationships of individual cells, with potential ramifications in emerging applications such as modeling tumor microenvironments for precision immunotherapy [41–43]. The standard approaches focus on instance segmentation by assigning unique identifiers to individual cells to facilitate downstream analysis [44–46]. Hover-net represents a significant advancement in addressing the limitations of semantic breadth and cell categorization within segmentation tasks, by incorporating cell classification into the segmentation process [47]. However, traditional methods typically rely on bounding box detection, and struggle with diverse cell morphologies and irregular shapes. Recent efforts aim to overcome these challenges by adopting more refined representations and accommodating the multi-resolution nature

of biological imaging [48–50]. Cell-ViT is a marquee example that leverages SAM’s encoder backbone to improve hierarchical representation, particularly for nucleus segmentation [51]. *BiomedParse* can contribute to this long line of exciting research work by enabling cell segmentation and identification in one fell swoop and enhancing generalizability through joint training on a diverse range of image modalities and cell types.

While *BiomedParse* has demonstrated promising potential for unifying biomedical image analysis, growth areas abound. First, although *BiomedParse* has demonstrated high accuracy (e.g., Dice scores) in identifying relevant pixels in an image for a given object type, by default it does not differentiate individual object instances and requires post processing to separate the instance masks, which is important in some applications such as cell counting. Second, while *BiomedParse* can already perform image analysis from text prompt alone, it currently does not support interactive dialogue with users in a conversational style like GPT-4. To address this, we plan to develop a conversational system that can better tailor to complex user needs. Finally, *BiomedParse* currently treats non-2D modalities such as CT and MRI by reducing them to 2D slices, thus failing to utilize the spatial and temporal information in the original modalities. In future work, we need to extend *BiomedParse* beyond 2D image slices to facilitate 3D segmentation, detection, and recognition.

## Methods

### Details of *BiomedParseData*

We created the first large-scale biomedical image parsing dataset *BiomedParseData*, where each image is associated with a collection of objects. Each object is annotated with the segmentation mask and a canonical semantic label specifying the object type from a biomedical object ontology. Additionally, each semantic label comes with a set of synonymous textual descriptions for model training. *BiomedParseData* was created by synthesizing 45 publicly available biomedical segmentation datasets across 9 imaging modalities, comprising 1.1 million images, 3.4 million image-mask-label triples, and 6.8 million image-mask-description triples (**Fig. 1b**). To ensure the quality of *BiomedParseData*, we imposed stringent inclusion criteria: each image had to be manually or semi-manually segmented at the pixel level, and a name was available for each segmented object from the dataset description. For 3D imaging modalities such as CT and MRI, we pre-processed each volume into in-plane 2D slices to be consistent with other modalities.

For model training and evaluation, we randomly split each original dataset into 80% training and 20% testing. Slices from each 3D volume always appear in the same split to prevent information leakage.

To harmonize natural-language variations in noisy object descriptions, we use GPT-4 to create a three-layer biomedical object ontology (**Fig. 1a**). The base layer comprises three broad semantic categories: organ, abnormality, histology. The next layer comprises 15 meta-object types (e.g., heart in organ and tumor in abnormality). The most fine-grained layer comprises 82 object types, such as left heart ventricle and enhancing tumor. Specifically, we first used GPT-4 to generate a preliminary hierarchical structure for biomedical image analysis and propose candidate names for individual object types, drawing from a wide range of tasks and textual descriptions across the source datasets. We then manually reviewed these candidates and mapped them to standardized OHDSI vocabularies using Athena [52]. We introduce *other* as a catch-all category. For future expansion, we expect that the first two layers are relatively stable, while our framework can easily incorporate new object types in the fine-grained layers, as well as additional datasets with segmentation and object labels.

To enhance the robustness of *BiomedParse* in handling diverse text prompts, we also used GPT-4 to generate synonymous textual descriptions for each semantic label, following other recent efforts in using GPT-4 for synthetic data generation [53, 54]. Specifically, we adopt a template normalization for each dataset, by formulating the unifying image analysis task as identifying “[OBJECT TYPE] in [ANATOMIC SITE] [MODALITY]”, such as “enhancing tumor in brain MRI” (**Supplementary Figure 2**). We then introduced linguistic diversity into these descriptions by using GPT-4 to generate variations in professional language (**Supplementary Figure 1**), as well as introducing synonymous variations for each component (**Supplementary Figure 3**). In each training epoch, we randomly sampled a description for each image-mask pair, ensuring *BiomedParse* to understand diverse text prompts.

### Details of *BiomedParse*

Existing image analysis methods often focus on segmentation alone. They typically expect spatial input prompts such as bounding box or scribble for the object to segment, and focus on learning spatial embedding such as bounding box coordinates [8, 9, 55]. In contrast, *BiomedParse* follows SEEM [56] and focuses on learning text prompt. Specifically, *BiomedParse* adopts a modular design, comprising an image encoder, a text encoder, a mask decoder, and a meta-object classifier. See **Fig. 1c**. The image and text encoders were initialized using Focal [17] and PubMedBERT [18], respectively.

The input to *BiomedParse* is an image and a text prompt, which are passed along to the image and text encoders, respectively. The text prompt specifies the object type for segmentation and detection. The mask decoder outputs a segmentation mask that has the same size as the original image, with a probability between

0 and 1 for each pixel, indicating how likely the pixel belongs to the designated object in the text prompt. The meta-object classifier includes input from the image and outputs the meta-object type to facilitate joint training of image encoder with object semantics.

## Implementation of competing methods

We compared *BiomedParse* to state-of-the-art segmentation models, SAM [8] and MedSAM [9]. We recognize the importance of precise bounding boxes as model input, so we evaluated competing methods in two settings: (i) employing gold-standard bounding boxes, and (ii) utilizing bounding boxes predicted by the state-of-the-art object detection model Grounding DINO [16] to provide bounding box prompts. For the first setting, we follow [9] by deriving bounding boxes from gold-standard masks, ensuring each box tightly encompassed the mask with a uniform margin of 10 pixels. In the second setting, we adhered to the inference pipeline of Grounding DINO where when presented with multiple bounding box predictions, we selected the one with the highest confidence score. This text-to-box-to-segmentation scheme follows the idea in [57]. To maintain uniformity across comparisons, all input images were resized to  $1024 \times 1024$  pixels. We use the same test split of *BiomedParseData* for evaluation across competing methods, and performance was quantified using the median Dice score on each task. We recognize that the train-test splits are different across the original evaluations of the competing methods, and the *BiomedParseData* test split could contain examples that were used to train other models. We note that the implementations for MedSAM, SAM, and Grounding DINO were used as-is for inference purposes without any fine-tuning. As for the task-specific nnU-Net models [26] and the DeepLabV3+ models [58], due to the unavailability of numerous task-finetuned models and the lack of explicit training details in existing literature, we relied on performance metrics reported in the MedSAM study [9].

## Detecting invalid textual description

*BiomedParse* by design can input any image and text prompt. However, a text prompt may be invalid, specifying an object that doesn't exist in the given image [54, 59]. For example, the request to identify and segment "left heart ventricle" in a dermoscopy image should be rejected by the model as invalid. It is critical to detect and reject invalid text prompt to preempt hallucinations [60].

In principle, the mask decoder should output low pixel probabilities for invalid text prompt. However, given the sheer number of pixels, some might get a relatively high output probability simply by chance, thus leading to erroneous object detection and segmentation results. To address this problem, we observe that while individual pixels might get noisily high probabilities, collectively their distribution would be rather different compared to pixels in valid objects. Consequently, we can estimate the distribution of its pixel probabilities from training data, and then estimate how likely the pixel probabilities in a test image are drawn from the same distribution.

Specifically, after *BiomedParse* was trained, for each object type, we computed the average object pixel probability for each training image containing objects of the given type, and fit a Beta distribution for all these probabilities. At test time, for a given image, we computed the average object pixel probability for the predicted object segments of the given object type, and compute the  $p$ -value using one-sample Kolmogorov-Smirnov (K-S) test [61]. Smaller  $p$ -value indicates that the predicted object segments are unlikely to be correct. To increase the robustness, in addition to pixel probability, we also consider the RGB values. In particular, for each color channel (R, G, B), we similarly fit a Beta distribution from the average value for valid objects in training, and compute the corresponding  $p$ -value for the predicted object segments in a test image. Overall, we treat these four tests as independent and use their product as the summary  $p$ -value.

In this way, we can obtain a summary  $p$ -value for any given pair of text prompt and image. To identify a summary  $p$ -value threshold for separating valid inputs from invalid ones, we created an invalid dataset by

sampling invalid object types for each image. We plot the distribution for both valid text prompts (for a given image) and invalid ones (**Fig. 2f**). For comparison against Grounding DINO, we use its confidence score given a text prompt and an image for invalid input detection.

### Attention map conditioned on the textual description

To visualize the shape of each segmentation object type, e.g. “hepatic vessel in CT”, we collected the predicted pixel probabilities for each object type and aggregated probabilities from all images. The pixel-level probability is derived from the top layer attention on the pixel. The attention map, reflecting the shape for a target  $t$ , is obtained in a four-step approach. First, we collected all *BiomedParse*-predicted pixel attention for target  $t$  as  $\rho_1, \dots, \rho_n \in [0, 1]^{H \times W}$  across  $n$  examples in the test set. Second, we initialized shape distribution for target  $t$  as  $\mathcal{M}_1^t = \rho_1$ . Third, for iteration  $i = 1, \dots, n - 1$ , we computed 2-D cross-correlation between  $\rho_{i+1}$  and  $\mathcal{M}_i^t$ , and shifted  $\rho_{i+1}$  to be aligned with  $\mathcal{M}_i^t$  at highest cross-correlation, and updated the ensemble distribution  $M_{i+1}^t = M_i^t + \tilde{\rho}_{i+1}$ , where  $\tilde{\rho}_{i+1}$  denotes the shifted attention matrix. Finally, the attention map for target  $t$  is normalized as  $M_n^t/n$ . For 3D segmentation targets such as CT and MRI, we first aggregated the predictions within one volume without shifting and then aligned the volume-aggregated masks using the above method.

### Details of experiments on irregular-shaped object detection

Medical image segmentation models like MedSAM require a bounding box as input. When the shape of the target is “irregular”, it is hard for the bounding box to precisely define the region of interest. To quantify the “regularity” of a target mask  $M$ , we define the following three metrics: **Box Ratio** measures the degree to which the target mask is similar to its tight bounding box:  $BoxRatio(M) = \frac{|M|}{|Box(M)|}$ , where  $Box(M)$  is the tight bounding box around mask  $M$ , and  $|\cdot|$  denotes the area measured in number of pixels. **Convex Ratio** measures how convex the target mask is and is defined as  $ConvexRatio(M) = \frac{|M|}{|ConvexHull(M)|}$ , where  $ConvexHull(M)$  is the convex hull of mask  $M$ . **Inverse rotational inertia** (IRI) measures how spread out the area of the target mask is. To begin with, the rotational inertia of  $M$  relative to its centroid  $c_M$  is  $RI(M) = \sum_{x \in M} \|x - c_M\|_2^2$ , where  $x$  is the coordinate of each pixel in the mask, and  $c_M$  is the coordinate of the centroid. To standardize the metric to be independent of the total mask area, we take the inverse of the rotational inertia and scale by the value of a round-shaped mask with the same area, representing the lowest rotational inertia achievable by any mask with the same area:  $IRI(M) = \frac{|M|^2}{2\pi \cdot RI(M)}$ . Under this definition, any mask has  $0 < IRI \leq 1$ , with any round-shaped mask having IRI equal to 1.

### Details of experiments on object recognition

We built a hierarchical structure putting all supported targets under one modality at one anatomic site. Given any image, e.g. abdominal CT, we traverse all the available targets  $t = 1, \dots, m$  under the branch that are exclusive to each other, and prompt the *BiomedParse* model sequentially to get  $m$  prediction of mask probabilities  $\rho^1, \dots, \rho^m$ . It is possible that the predicted masks can overlap with each other. The challenges then are how to select the right set of targets in the specific image and how to determine the right mask regions for the selected targets to avoid overlapping. We used a two-stage approach for object recognition, including a target selection stage and a mask aggregation stage. In the target selection stage, we first calculate the original mask area for each target  $t$  as  $A^t$ . Then, we iterate through the pixels. For each pixel  $(i, j)$ , we rank the targets that have pixel probability  $\rho_{ij}^t > 0.5$ . The target assigned to pixel  $(i, j)$  is  $T_{ij} = \text{argmax } \rho_{ij}^t$ . After this round of pixel assigning, the final area for each target  $t$  is  $\tilde{A}^t = \sum_{i,j} \mathbf{1}_{T_{ij}=t}$ . The targets with final area  $\tilde{A}^t > \lambda A^t$  are the selected targets, with  $\lambda$  being the user-specified threshold. In

the mask aggregation stage, we discard all unselected target masks completely, and then iterate through the pixels again. For each pixel, the most probable target  $t$  with  $\rho_{ij}^t > 0.5$  is assigned. The pixels with predicted probabilities  $\rho_{ij}^t \leq 0.5$  for all selected targets are left blank.

For the baseline method using Grounding DINO with SAM and MedSAM, we first prompted Grounding DINO with the set of targets to retrieve a collection of bounding boxes with confidence scores. Then we implemented non-maximum suppression [62–64] to select a subset of identified targets in the scene, minimizing the overlapping between the targets. To get the segmentation masks for these identified targets, we further prompted SAM and MedSAM with the bounding boxes to retrieve the corresponding predictions.

## Data availability

We will provide access to *BiomedParseData* or scripts to reproduce *BiomedParseData* from the original datasets, upon publication of this manuscript.

## Code availability

*BiomedParse* will be made fully available upon publication, including the model weights and relevant source code for pre-training, fine-tuning, and inference. We will also provide detailed methods and implementation steps to facilitate independent replication.

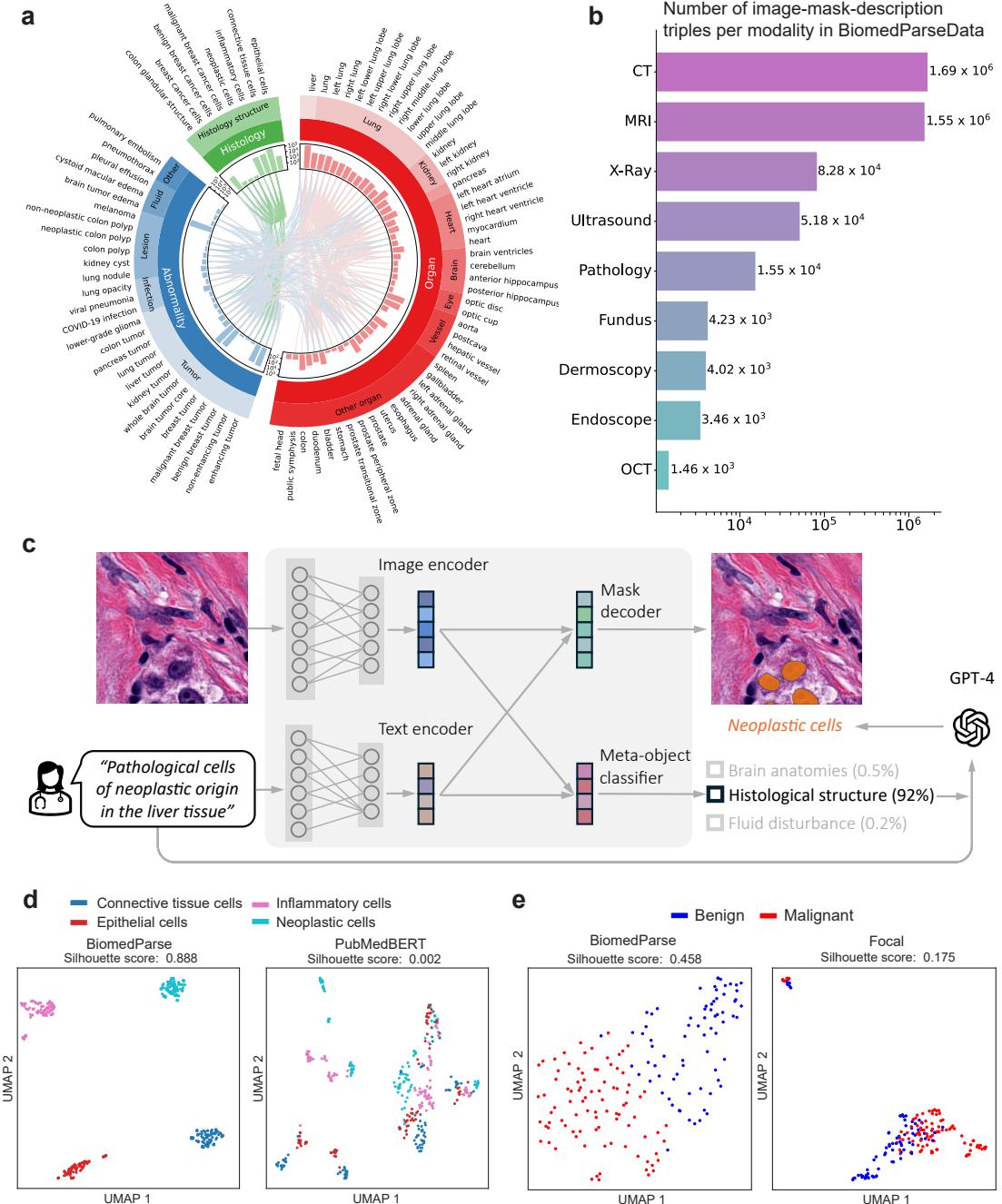
## References

- [1] Royer, L. A. The future of bioimage analysis: a dialog between mind and machine. *Nature Methods* **20**, 951–952 (2023).
- [2] Li, X., Zhang, Y., Wu, J. & Dai, Q. Challenges and opportunities in bioimage analysis. *Nature Methods* **20**, 958–961 (2023).
- [3] Wang, R. *et al.* Medical image segmentation using deep learning: A survey. *IET Image Processing* **16**, 1243–1267 (2022).
- [4] Salpea, N., Tzouveli, P. & Kollias, D. Medical image segmentation: A review of modern architectures. In *European Conference on Computer Vision*, 691–708 (Springer, 2022).
- [5] Ribli, D., Horváth, A., Unger, Z., Pollner, P. & Csabai, I. Detecting and classifying lesions in mammograms with deep learning. *Scientific reports* **8**, 4165 (2018).
- [6] Ma, W., Lu, J. & Wu, H. Cellcano: supervised cell type identification for single cell atac-seq data. *Nature Communications* **14**, 1864 (2023).
- [7] Jiang, H. *et al.* A review of deep learning-based multiple-lesion recognition from medical images: classification, detection and segmentation. *Computers in Biology and Medicine* **157**, 106726 (2023).
- [8] Kirillov, A. *et al.* Segment anything. *arXiv preprint arXiv:2304.02643* (2023).
- [9] Ma, J. *et al.* Segment anything in medical images. *Nature Communications* **15**, 654 (2024).
- [10] Tu, Z., Chen, X., Yuille, A. L. & Zhu, S.-C. Image parsing: Unifying segmentation, detection, and recognition. *International Journal of computer vision* **63**, 113–140 (2005).
- [11] Tighe, J. & Lazebnik, S. Superparsing: scalable nonparametric image parsing with superpixels. *International Journal of Computer Vision* **101**, 329–349 (2013).
- [12] Zhou, S. K. *Medical image recognition, segmentation and parsing: machine learning and multiple object approaches* (Academic Press, 2015).
- [13] Gamper, J. *et al.* Pannuke dataset extension, insights and baselines. *arXiv preprint arXiv:2003.10778* (2020).
- [14] Ji, Y. *et al.* Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation. *arXiv preprint arXiv:2206.08023* (2022).
- [15] Bernard, O. *et al.* Deep learning techniques for automatic MRI cardiac multi-structures segmentation and diagnosis: Is the problem solved? *IEEE Transactions on Medical Imaging* **37**, 2514–2525 (2018).
- [16] Liu, S. *et al.* Grounding DINO: Marrying DINO with grounded pre-training for open-set object detection (2023). [2303.05499](https://arxiv.org/abs/2303.05499).
- [17] Yang, J., Li, C., Dai, X. & Gao, J. Focal modulation networks. *Advances in Neural Information Processing Systems* **35**, 4203–4217 (2022).
- [18] Gu, Y. *et al.* Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)* **3**, 1–23 (2021).

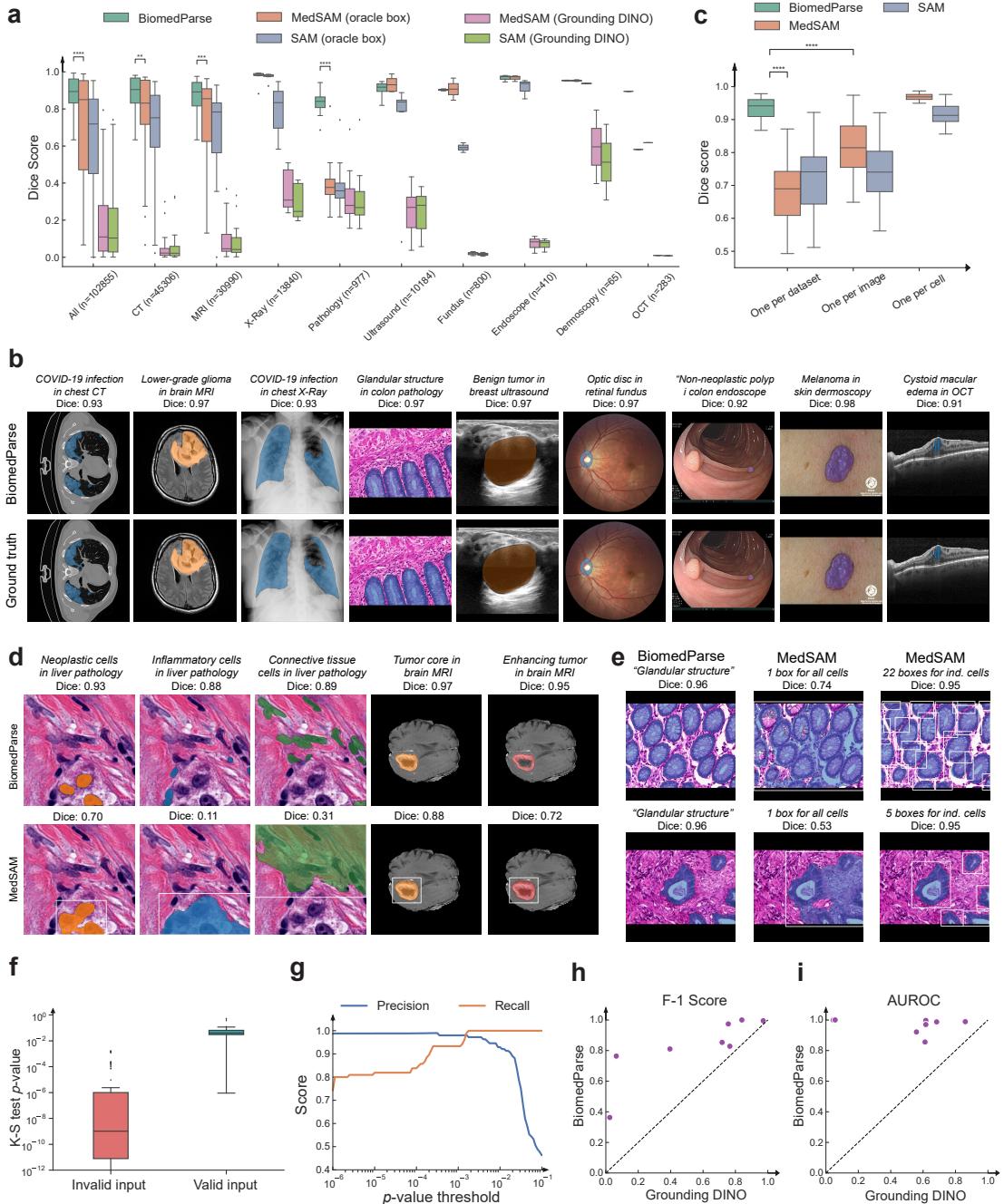
- [19] Sirinukunwattana, K., Snead, D. R. J. & Rajpoot, N. M. A stochastic polygons model for glandular structures in colon histology images. *IEEE Transactions on Medical Imaging* **34**, 2366–2378 (2015).
- [20] Ronneberger, O., Fischer, P. & Brox, T. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III* 18, 234–241 (Springer, 2015).
- [21] Çiçek, Ö., Abdulkadir, A., Lienkamp, S. S., Brox, T. & Ronneberger, O. 3D U-Net: learning dense volumetric segmentation from sparse annotation. In *International conference on medical image computing and computer-assisted intervention*, 424–432 (Springer, 2016).
- [22] Milletari, F., Navab, N. & Ahmadi, S.-A. V-Net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, 565–571 (IEEE, 2016).
- [23] Li, X. *et al.* H-DenseUNet: hybrid densely connected unet for liver and tumor segmentation from CT volumes. *IEEE transactions on medical imaging* **37**, 2663–2674 (2018).
- [24] Zhou, Z., Siddiquee, M. M. R., Tajbakhsh, N. & Liang, J. UNet++: Redesigning skip connections to exploit multiscale features in image segmentation. *IEEE transactions on medical imaging* **39**, 1856–1867 (2019).
- [25] Myronenko, A. 3D MRI brain tumor segmentation using autoencoder regularization. In *International MICCAI Brainlesion Workshop*, 311–320 (Springer, 2018).
- [26] Isensee, F., Jaeger, P. F., Kohl, S. A., Petersen, J. & Maier-Hein, K. H. nnU-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods* **18**, 203–211 (2021).
- [27] Lee, H. H., Bao, S., Huo, Y. & Landman, B. A. 3D UX-Net: A large kernel volumetric convnet modernizing hierarchical transformer for medical image segmentation. *arXiv preprint arXiv:2209.15076* (2022).
- [28] Lee, H. H. *et al.* Scaling up 3D kernels with bayesian frequency re-parameterization for medical image segmentation. *arXiv preprint arXiv:2303.05785* (2023).
- [29] Lee, H. H. *et al.* DeformUX-Net: Exploring a 3D foundation backbone for medical image segmentation with depthwise deformable convolution. *arXiv preprint arXiv:2310.00199* (2023).
- [30] Chen, J. *et al.* TransUNet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306* (2021).
- [31] Xu, G., Zhang, X., He, X. & Wu, X. LeViT-UNet: Make faster encoders with transformer for medical image segmentation. In *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*, 42–53 (Springer, 2023).
- [32] Xie, Y., Zhang, J., Shen, C. & Xia, Y. Cotr: Efficiently bridging CNN and transformer for 3D medical image segmentation. In *International conference on medical image computing and computer-assisted intervention*, 171–180 (Springer, 2021).
- [33] Wang, W. *et al.* TransBTS: Multimodal brain tumor segmentation using transformer. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 109–119 (Springer, 2021).

- [34] Hatamizadeh, A. *et al.* UNETR: Transformers for 3D medical image segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 574–584 (2022).
- [35] Hatamizadeh, A. *et al.* Swin UNETR: Swin transformers for semantic segmentation of brain tumors in MRI images. In *International MICCAI Brain lesion Workshop*, 272–284 (Springer, 2022).
- [36] Zhou, H.-Y. *et al.* nnFormer: Interleaved transformer for volumetric segmentation. *arXiv preprint arXiv:2109.03201* (2021).
- [37] Cao, H. *et al.* Swin-UNet: UNet-like pure transformer for medical image segmentation. *arXiv preprint arXiv:2105.05537* (2021).
- [38] Ren, S., He, K., Girshick, R. & Sun, J. Faster r-cnn: Towards real-time object detection with region proposal networks (2016). [1506.01497](https://arxiv.org/abs/1506.01497).
- [39] Bochkovskiy, A., Wang, C.-Y. & Liao, H.-Y. M. Yolov4: Optimal speed and accuracy of object detection (2020). [2004.10934](https://arxiv.org/abs/2004.10934).
- [40] Litjens, G. *et al.* A survey on deep learning in medical image analysis. *Medical Image Analysis* **42**, 60–88 (2017). URL <http://dx.doi.org/10.1016/j.media.2017.07.005>.
- [41] Stringer, C., Wang, T., Michaelos, M. & Pachitariu, M. Cellpose: a generalist algorithm for cellular segmentation. *Nature methods* **18**, 100–106 (2021).
- [42] Greenwald, N. F. *et al.* Whole-cell segmentation of tissue images with human-level performance using large-scale data annotation and deep learning. *Nature biotechnology* **40**, 555–565 (2022).
- [43] Ma, J. & Wang, B. Towards foundation models of biological image segmentation. *Nature Methods* **20**, 953–955 (2023).
- [44] Girshick, R. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, 1440–1448 (2015).
- [45] He, K., Gkioxari, G., Dollár, P. & Girshick, R. Mask R-CNN. In *Proceedings of the IEEE international conference on computer vision*, 2961–2969 (2017).
- [46] Schmidt, U., Weigert, M., Broaddus, C. & Myers, G. Cell detection with star-convex polygons. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part II* **11**, 265–273 (Springer, 2018).
- [47] Graham, S. *et al.* Hover-Net: Simultaneous segmentation and classification of nuclei in multi-tissue histology images. *Medical image analysis* **58**, 101563 (2019).
- [48] Yang, H. *et al.* CircleNet: Anchor-free glomerulus detection with circle representation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part IV* **23**, 35–44 (Springer, 2020).
- [49] Nguyen, E. H. *et al.* CircleSnake: Instance segmentation with circle representation. In *International Workshop on Machine Learning in Medical Imaging*, 298–306 (Springer, 2022).
- [50] Ilyas, T. *et al.* Tsfd-net: Tissue specific feature distillation network for nuclei segmentation and classification. *Neural Networks* **151**, 1–15 (2022).

- [51] Hörst, F. *et al.* Cellvit: Vision transformers for precise cell segmentation and classification. *arXiv preprint arXiv:2306.15350* (2023).
- [52] OHDSI. Athena standardized vocabularies. <https://www.ohdsi.org/analytic-tools/athena-standardized-vocabularies/> (n.d.). Accessed: 2022-01-17.
- [53] Gu, Y. *et al.* BiomedJourney: Counterfactual biomedical image generation by instruction-learning from multimodal patient journeys. *arXiv preprint arXiv:2310.10765* (2023).
- [54] Li, C. *et al.* LLaVA-Med: Training a large language-and-vision assistant for biomedicine in one day. *Advances in Neural Information Processing Systems* **36** (2024).
- [55] Wong, H. E., Rakic, M., Guttag, J. & Dalca, A. V. Scribbleprompt: Fast and flexible interactive segmentation for any medical image. *arXiv preprint arXiv:2312.07381* (2023).
- [56] Zou, X. *et al.* Segment everything everywhere all at once. *Advances in Neural Information Processing Systems* **36** (2024).
- [57] Ren, T. *et al.* Grounded SAM: Assembling open-world models for diverse visual tasks. *arXiv preprint arXiv:2401.14159* (2024).
- [58] Chen, L.-C., Zhu, Y., Papandreou, G., Schroff, F. & Adam, H. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, 801–818 (2018).
- [59] Lee, P., Goldberg, C. & Kohane, I. *The AI revolution in medicine: GPT-4 and beyond* (Pearson, 2023).
- [60] Achiam, J. *et al.* GPT-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [61] Massey Jr, F. J. The Kolmogorov-Smirnov test for goodness of fit. *Journal of the American statistical Association* **46**, 68–78 (1951).
- [62] Canny, J. A computational approach to edge detection. *IEEE Transactions on pattern analysis and machine intelligence* 679–698 (1986).
- [63] Viola, P. & Jones, M. Rapid object detection using a boosted cascade of simple features. In *Proceedings of the 2001 IEEE computer society conference on computer vision and pattern recognition. CVPR 2001*, vol. 1, I–I (Ieee, 2001).
- [64] Girshick, R., Donahue, J., Darrell, T. & Malik, J. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 580–587 (2014).

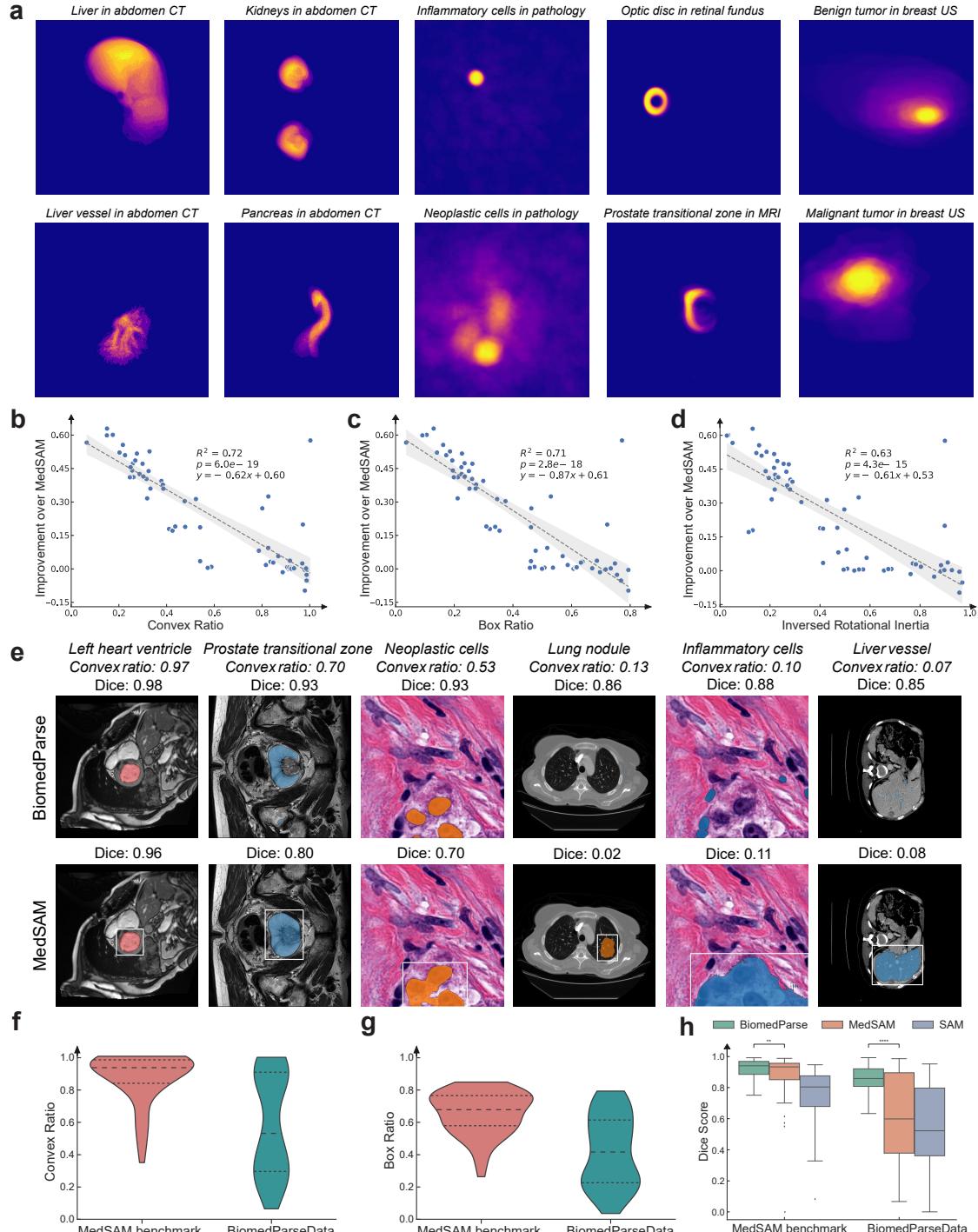


**Figure 1: Overview of *BiomedParse* and *BiomedParseData*.** **a**, The GPT-4 constructed ontology showing a hierarchy of object types that are used to unify semantic concepts across datasets. Bar plots showing the number of images containing that object type. Lines between bars showing the object type similarity in the text embedding space. **b**, Bar plot showing the number of image-mask-description triples for each modality in *BiomedParseData*. **c**, Flowchart of *BiomedParse*. *BiomedParse* takes an image and a text prompt as input and then outputs the segmentation masks for the objects specified in the prompt. Image-specific manual interaction such as bounding box or clicks is not required in our framework. To facilitate semantic learning for the image encoder, *BiomedParse* also incorporates a learning objective to classify the meta-object type. For evaluation, GPT-4 is used to resolve text prompt into object types using the object ontology, which also uses the meta-object type output from *BiomedParse* to narrow down candidate semantic labels. **d**, UMAP plots contrasting the text-embeddings for different cell types derived from *BiomedParse* text encoder (left) and PubMedBERT (right). **e**, UMAP plots contrasting the image embeddings for different cell types derived from *BiomedParse* image encoder (left) and Focal (right).

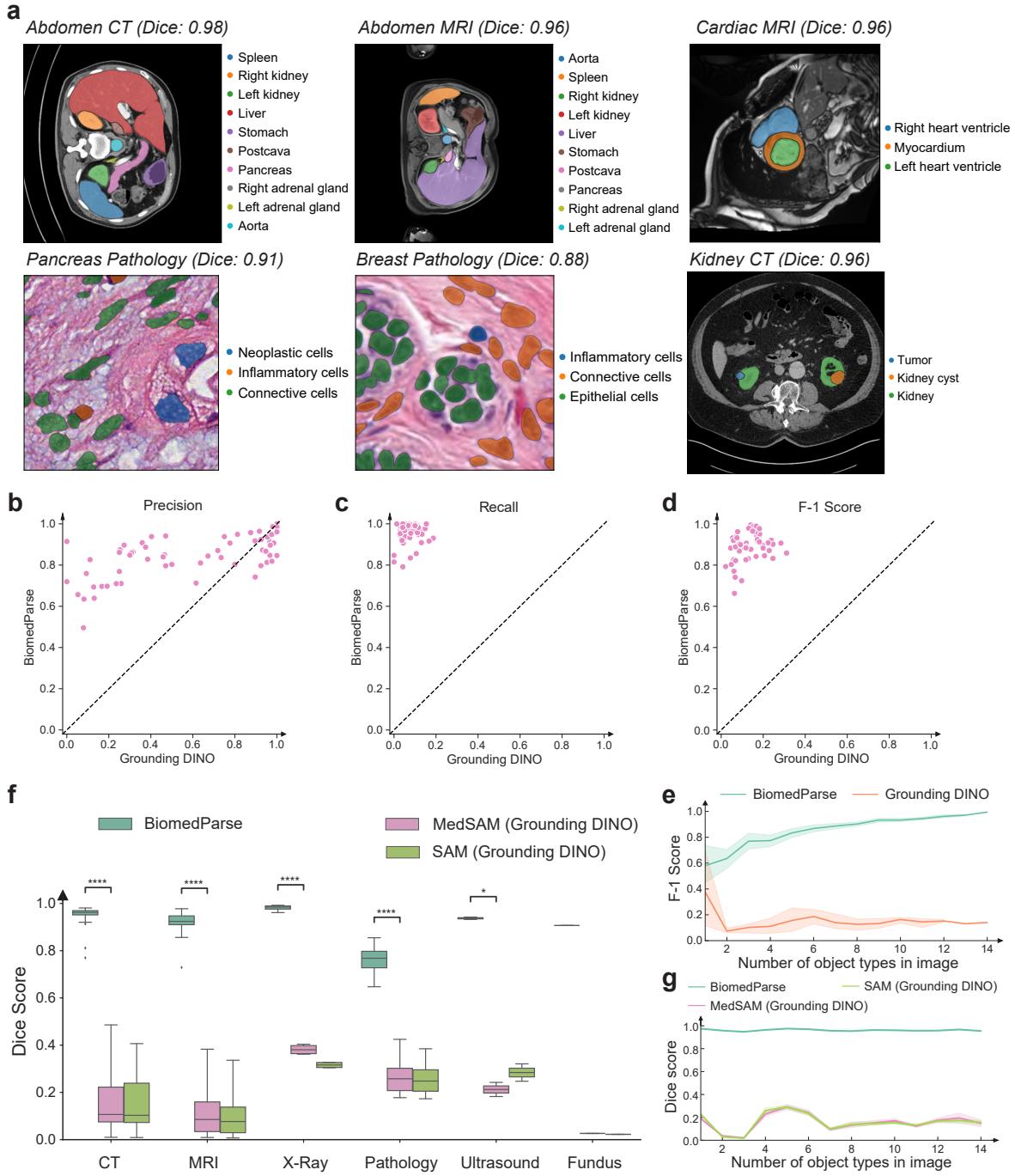


**Figure 2: Comparison on large-scale biomedical image segmentation datasets.** **a**, Bar plot comparing the Dice score between our method and competing methods on 102,855 test instances (image-mask-label triples) across 9 modalities. MedSAM and SAM require bounding box as input. We consider two settings: oracle bounding box (minimum bounding box covering the gold mask); bounding boxes generated from the text prompt by Grounding DINO, a state-of-the-art text-based grounding model.  $n$  denotes the number of test instances in the corresponding modality. \* indicates the significance level at which *BiomedParse* outperforms the best-competing method, with Wilcoxon test  $p$ -value  $< 1 \times 10^{-2}$  for \*\*\*,  $p$ -value  $< 1 \times 10^{-3}$  for \*\*\*\*,  $p$ -value  $< 1 \times 10^{-4}$  for \*\*\*\*\*. **b**, Nine examples comparing the segmentation results by *BiomedParse* and the ground truth, using just the text prompt at the top. **c**, Bar plot comparing the Dice score between our method and competing methods on a cell segmentation test set with 42 images. *BiomedParse* requires only a single user operation (the text prompt “Glandular structure in colon pathology”). By contrast, to get competitive results, MedSAM/SAM require 430 operations (one bounding box per an individual cell). **d**, Segmentation examples for liver pathology. **e**, Segmentation examples for brain MRI. **f**, KS test p-value distribution. **g**, Precision and Recall vs p-value threshold. **h**, F-1 Score vs Grounding DINO. **i**, AUROC vs Grounding DINO.

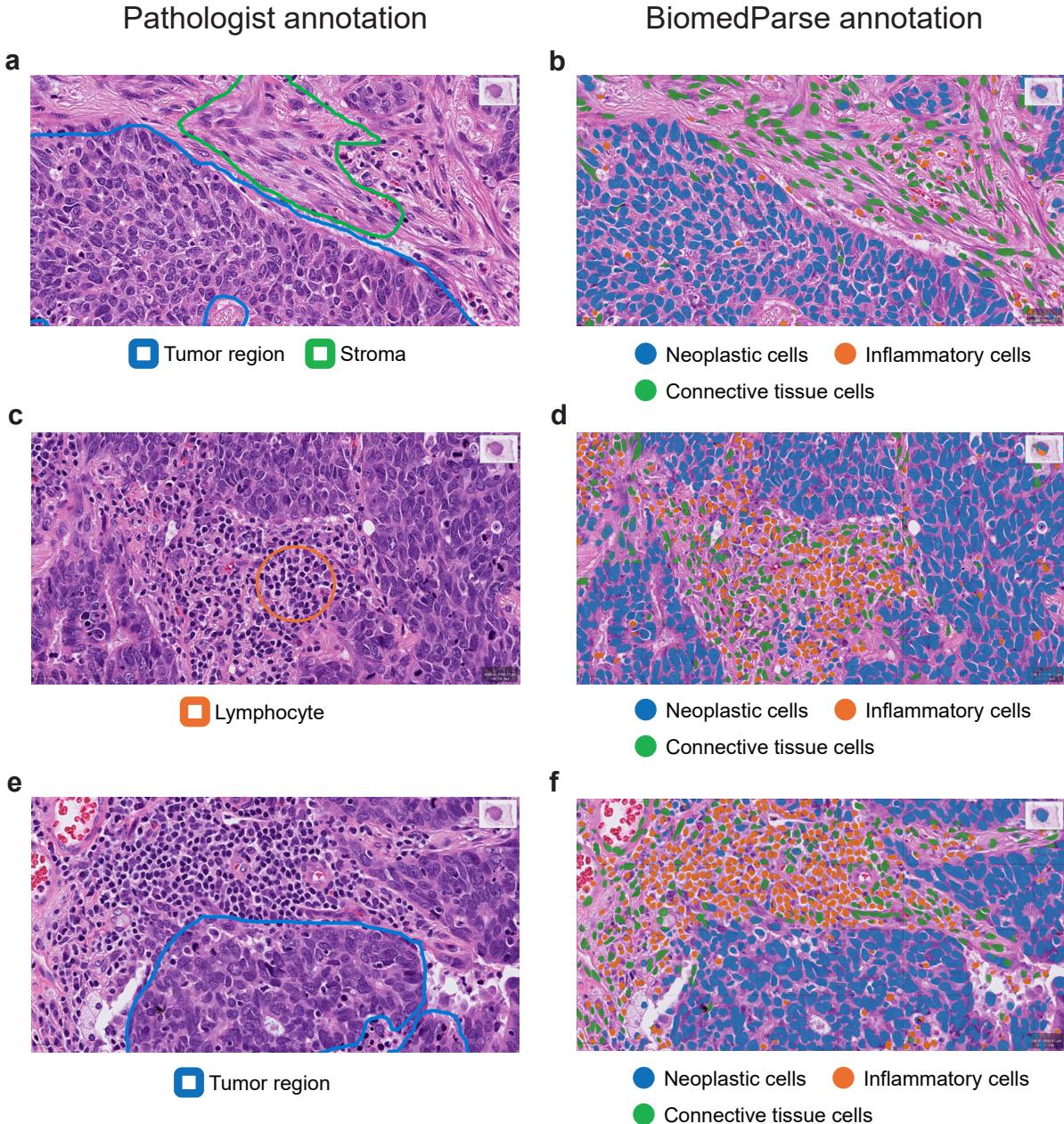
**d**, Five examples contrasting the segmentation results by *BiomedParse* and MedSAM, along with text prompts used by *BiomedParse* and bounding boxes used by MedSAM. **e**, Comparison between *BiomedParse* and MedSAM on a benign tumor image (top) and a malignant tumor image (bottom). The improvement of *BiomedParse* over MedSAM is even more pronounced on abnormal cells with irregular shapes. **f**, Bar plot comparing the K-S test  $p$ -values between valid text prompt and invalid text prompt. *BiomedParse* learns to reject invalid text prompts describing object types not present in the image (small  $p$ -value). **g**, Plot showing the precision and recall of our method on detecting invalid text prompts across different K-S test  $p$ -value cutoff. **h,i**, Scatter plots comparing the AUROC (**h**) and F-1 (**i**) between *BiomedParse* and Grounding DINO on detecting invalid descriptions.



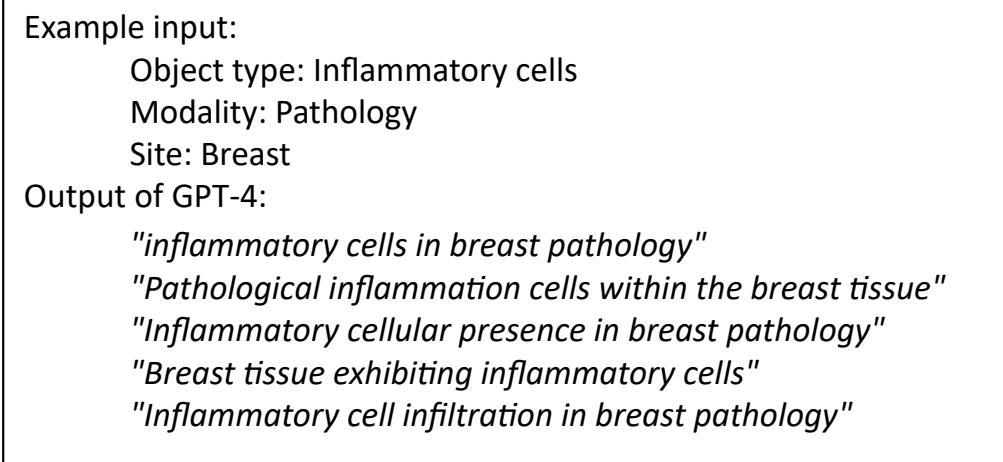
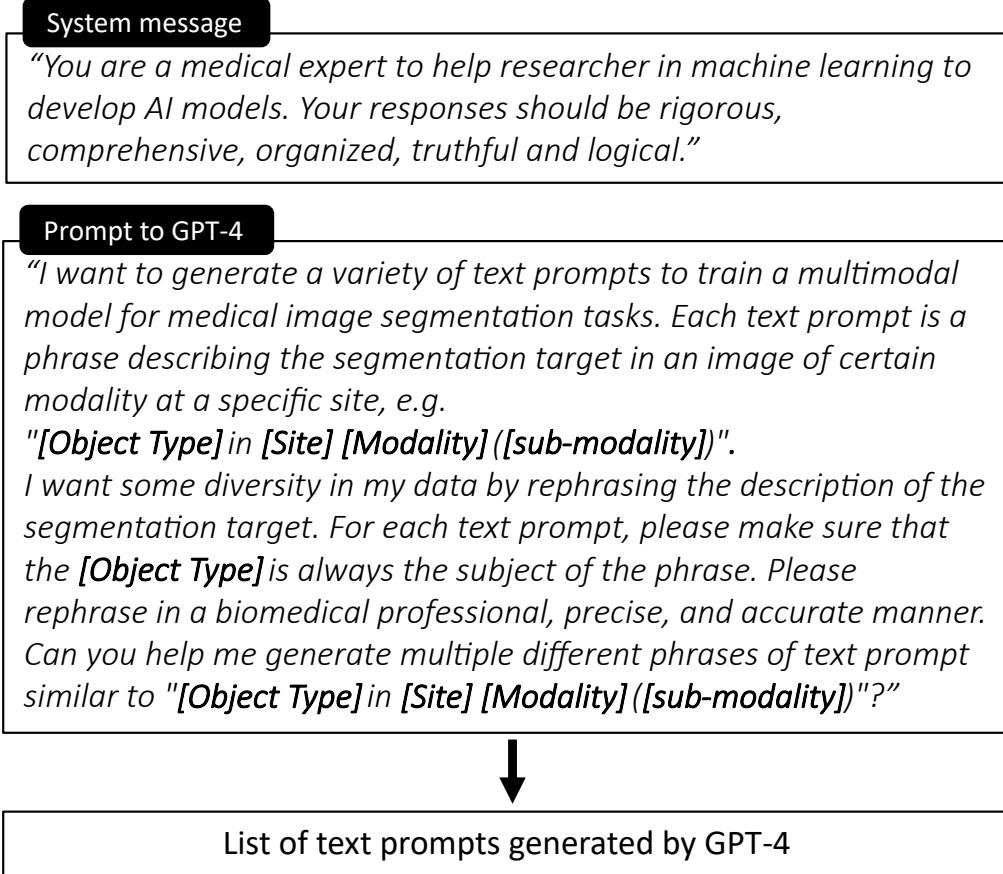
**Figure 3: Evaluation on detecting irregular-shaped objects.** **a**, Attention maps of text prompts for irregular-shaped objects, suggesting that *BiomedParse* learns rather faithful representation of their typical shapes. **b-d**, Scatter plots comparing the improvement in Dice score for *BiomedParse* over MedSAM with shape regularity in terms of convex ratio (**b**), box ratio (**c**), and inversed rotational inertia (**d**). Smaller number in x-axis means higher irregularity in average. Each dot is an object type. **e**, Six examples contrasting *BiomedParse* and MedSAM on detecting irregular-shaped objects. Plots are ordered from the least irregular one (left) to the most irregular one (right). **f,g** Comparison between *BiomedParseData* and the benchmark dataset used by MedSAM in terms of convex ratio (**f**) and box ratio (**g**). *BiomedParseData* is a more faithful representation of real-world challenges in terms of irregular-shaped objects. **h**, Bar plots comparing *BiomedParse* and competing approaches on *BiomedParseData* and the benchmark dataset used by MedSAM. *BiomedParse* has a larger improvement on *BiomedParseData*, which contains more diverse images and more irregular-shaped objects.



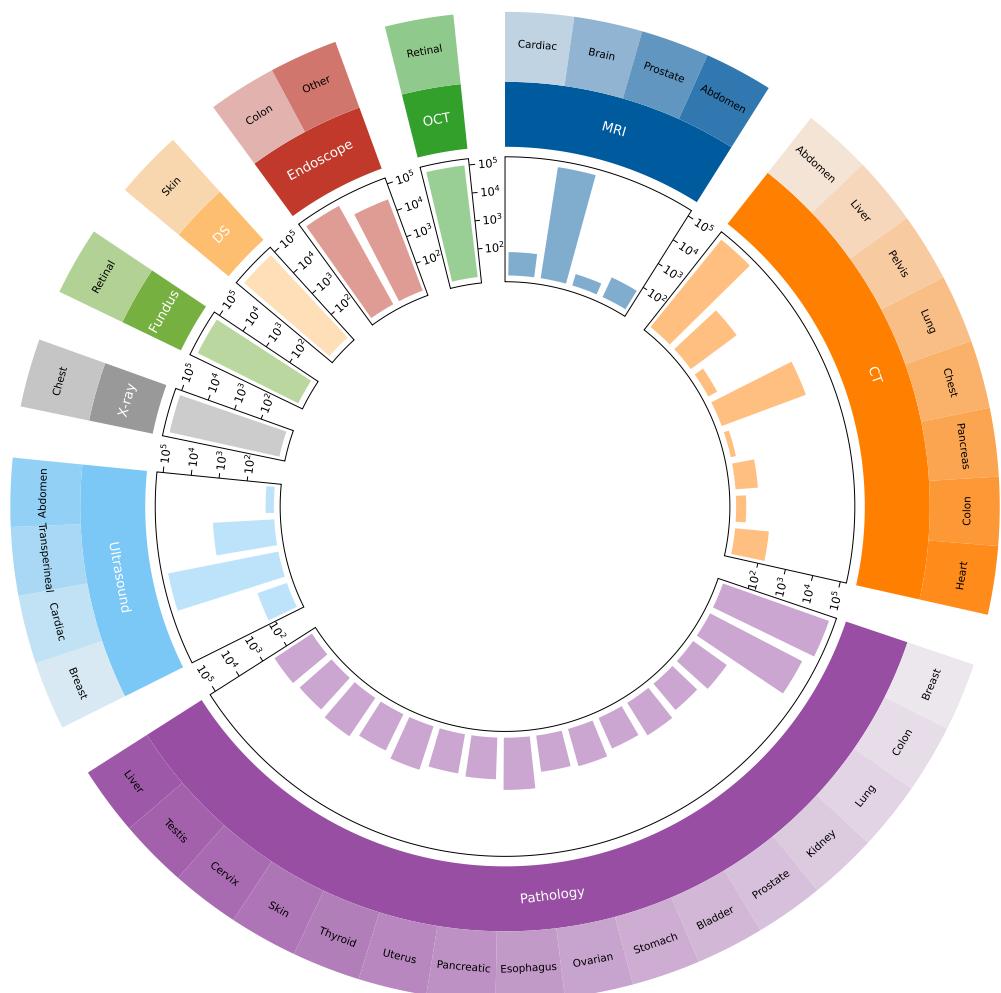
**Figure 4: Evaluation on object recognition.** **a**, Six examples showing the results of object recognition by our method. Object recognition identifies and segments all objects in an image without requiring any user-provided input prompt. **b-d**, Scatter plots comparing the F1 (**b**), Precision (**c**), and Recall (**d**) scores between *BiomedParse* and Grounding DINO on identifying objects presented in the image. **e**, Comparison between *BiomedParse* and Grounding DINO on object identification in terms of median F-1 score across different numbers of objects in the image. **f**, Bar plot comparing *BiomedParse* and MedSAM/SAM (using bounding boxes generated by Grounding DINO) on end-to-end object recognition (including segmentation) in relation to various modalities. **g**, Comparison between *BiomedParse* and MedSAM/SAM (using bounding boxes generated by Grounding DINO) on end-to-end object recognition (including segmentation) in relation to numbers of distinct objects in the image.



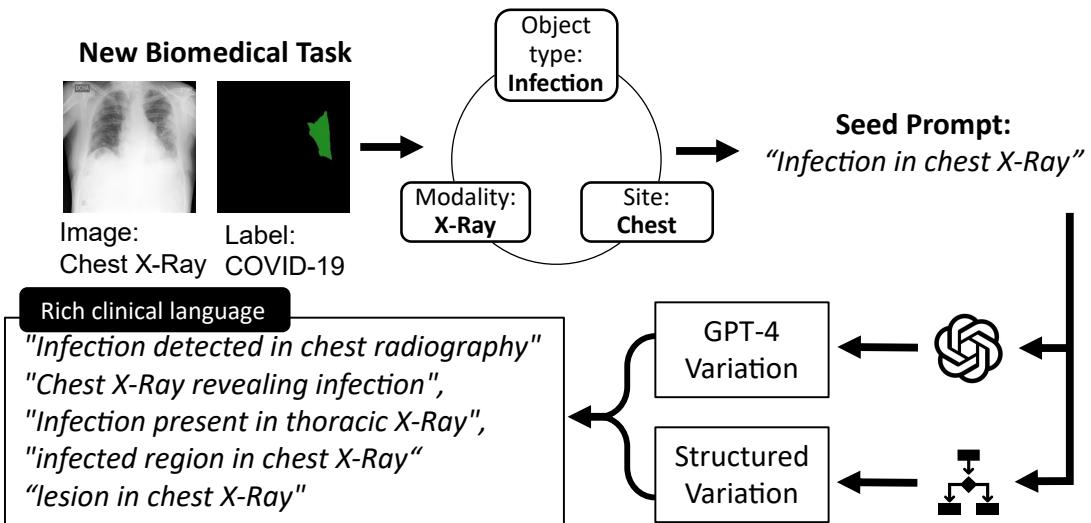
**Figure 5: Evaluation of *BiomedParse* on real-world cell segmentation examples.** **a-f**, De-identified pathology images from Providence Health System are used to compare the pathologist annotations (**a,c,e**) and the annotations from *BiomedParse* (**b,d,f**). We show the exact pathologist outputs, including object names (e.g., lymphocyte, stroma) and object locations, as well as the exact outputs by *BiomedParse*. *BiomedParse* does not need any user-provided inputs and can identify and segment cells of any types included in the ontology.



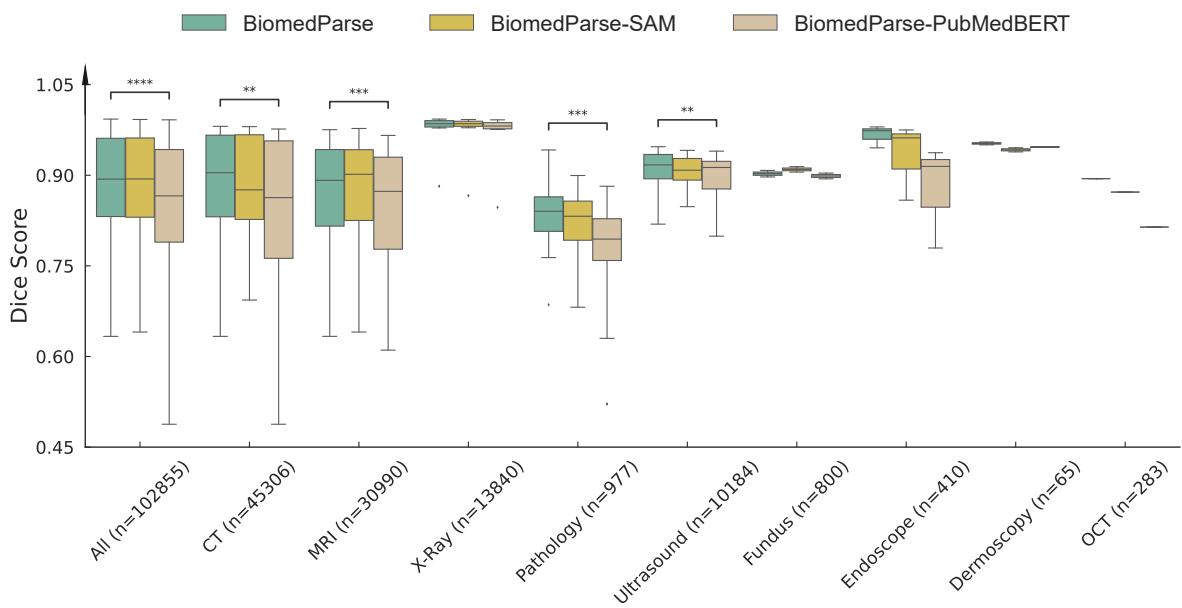
**Supplementary Figure 1:** GPT-4 prompt that is used to generate diverse descriptions for a given image according to its object type, image modality, and anatomic site.



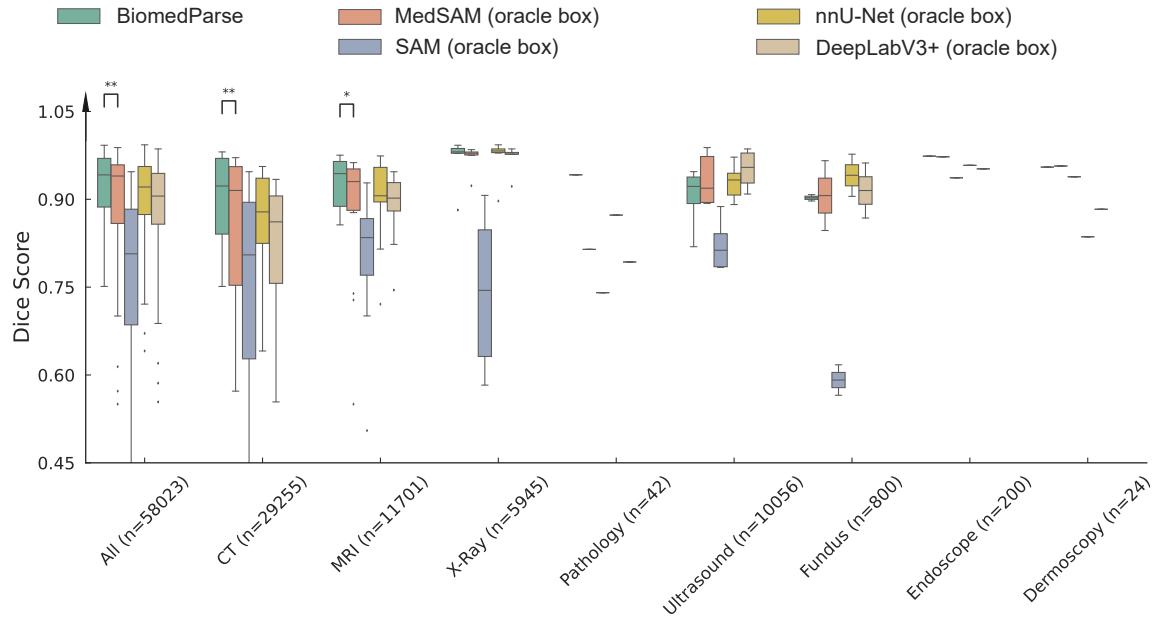
**Supplementary Figure 2:** Number of images in each of the 25 anatomic sites from 9 modalities. One anatomic site could present in multiple modalities.



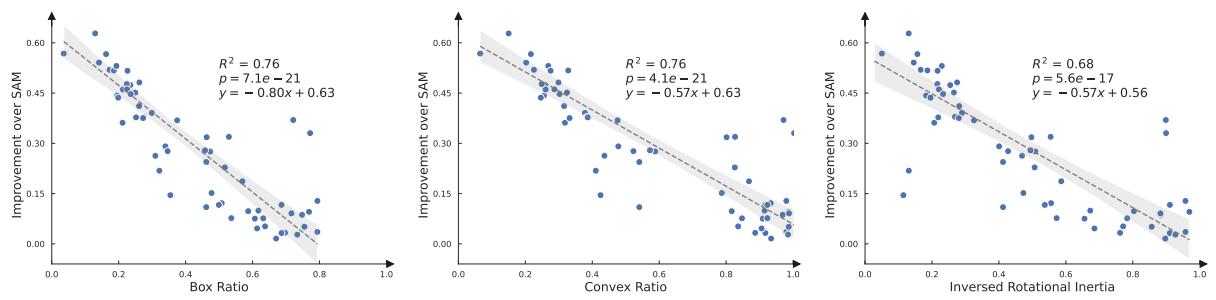
**Supplementary Figure 3:** Generating textual description for each object in each image. Object type, modality, and site are extracted from the metadata or the data description. We utilized both GPT-4 and structured biomedical concepts to generate rich variations of clinical language, increasing the robustness of *BiomedParse* to user-provided text.



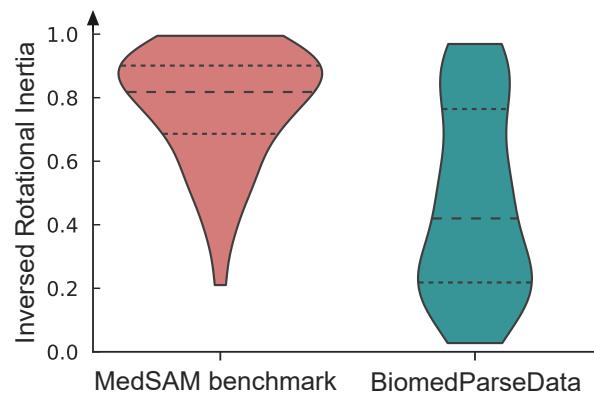
**Supplementary Figure 4:** Ablation studies comparing the performance of *BiomedParse* and two variants. *BiomedParse-SAM* stands for using SAM to initialize the image encoder. *BiomedParse-PubmedBERT* stands for using the frozen PubmedBERT [18] as the text encoder. *n* denotes the number of images in the corresponding modality. \* indicates the significance level at which *BiomedParse* outperforms the best-competing method, with Wilcoxon test  $p$ -value  $< 1 \times 10^{-2}$  for \*\*,  $p$ -value  $< 1 \times 10^{-3}$  for \*\*\*,  $p$ -value  $< 1 \times 10^{-4}$  for \*\*\*\*.



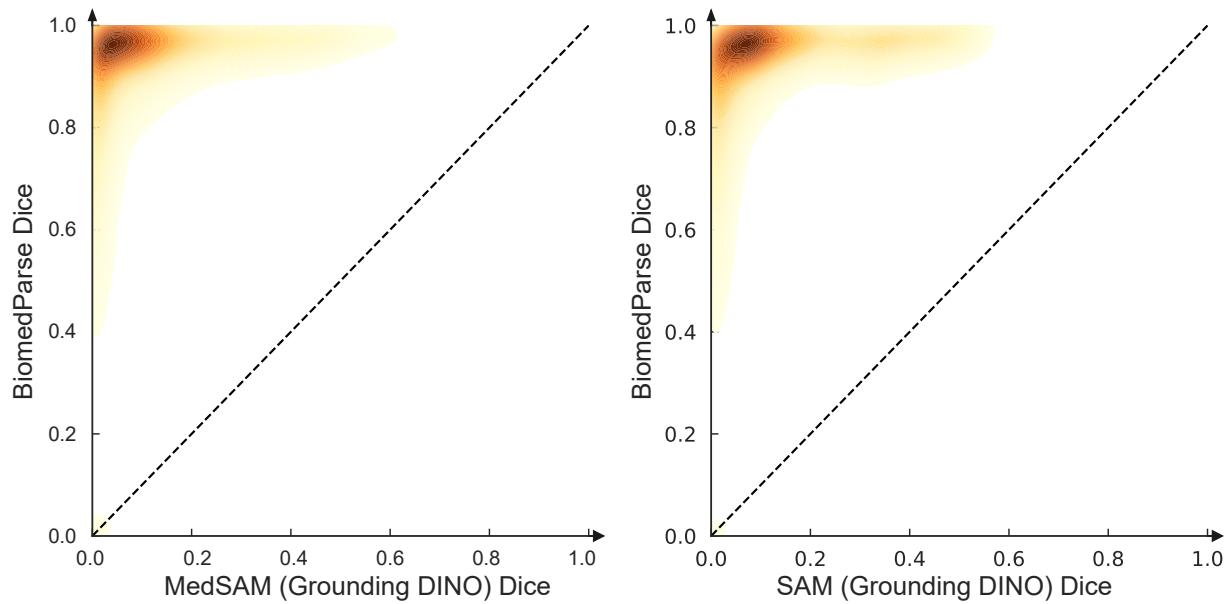
**Supplementary Figure 5:** Comparison between *BiomedParse* and competing methods on the MedSAM benchmark. We evaluated MedSAM and SAM using the ground truth bounding box for the segmentation. For nnU-Net and DeepLabV3+, we reported the evaluation reported by MedSAM [9]. Results are shown by imaging modality, with statistical significance comparison between *BiomedParse* and best-competing method MedSAM. \* indicates the significance level at which *BiomedParse* outperforms the best-competing method, with Wilcoxon test  $p\text{-value} < 5 \times 10^{-2}$  for \*,  $p\text{-value} < 1 \times 10^{-3}$  for \*\*.



**Supplementary Figure 6:** Scatter plots comparing the improvement of *BiomedParse* over SAM with shape irregularity in terms of box ratio (left), convex ratio (middle), and inversed rotational inertia (right). Each dot represents the median statistics over one object type in our segmentation ontology.



**Supplementary Figure 7:** Violin plot comparing the inversed rotational inertia between MedSAM benchmark data and *BiomedParseData*. A higher inversed rotational inertia indicates less irregularity.



**Supplementary Figure 8: Evaluation on object recognition.** **a,b,** Density plots comparing the performance on object recognition between *BiomedParse* and MedSAM (Grounding DINO) **a**, and between *BiomedParse* and SAM (Grounding DINO) **b**.