# A. Data Pre-processing

We make use of the publicly available anatomy-centered annotations provided by Chest ImaGemone [9, 10] to obtain masking maps of anatomical regions ($\{\mathcal{M}_n^{(\text{Abn})}\}_{n=1}^N, \mathcal{M}^{(\text{Chest})}$) and derive the hierarchical diagnostic description $\{H_m\}_{m=0}^M$ for each report in the MIMIC CXR dataset that is used for our performance evaluation.

The Chest ImaGenome dataset (https://physionet.org/content/chest-imagenome/1.0.0) contains scene graphs of each frontal chest X-ray image in the MIMIC CXR dataset. Each scene graph contains the bounding box coordinates of 29 anatomical regions presented in the image. The bounding box is associated with several sentences of the report if these sentences is about the clinical findings observed on that anatomical part. The abnormality terms mentioned in each sentence is also annotated.

## A.1. Constructing Masking Maps of Anatomical Regions

To obtain the masking maps, we first collect the coordinates of the pre-defined $N$ anatomical regions, denoted as $\{(\text{x}_n, \text{y}_n, \text{w}_n, \text{h}_n)\}_{n=1}^N$.

Consequently, given an X-ray image of the size of $W \times H$, we propose to obtain a masking map $\mathcal{M}_n^{(\text{Ana})} \in \mathbb{R}^{\mathcal{W} \times \mathcal{H}}$ for each anatomical parts with $\{0, 1\}$ value indicating whether the bounding box $(\text{x}_n, \text{y}_n, \text{w}_n, \text{h}_n)$ of the detected anatomical part falls on the corresponding feature region of the size of $\mathcal{H} \times \mathcal{W}$. The implementation is described as follows:

i) $(\text{x}_n, \text{y}_n, \text{w}_n, \text{h}_n)$ the coordinates of each anatomical part are re-scaled from $W \times H$ to $\mathcal{W} \times \mathcal{H}$, denoted as $(\text{x}_{1|n}, \text{x}_{2|n}, \text{y}_{1|n}, \text{y}_{2|n})$;

ii) An all-zero matrix with size $\mathcal{M}_n^{(\text{abn})} \in \mathbb{R}^{\mathcal{W} \times \mathcal{H}}$ is initialized for each anatomical part, and the values of the matrix are set to 1, denoted as $\mathcal{M}_n^{(\text{Ana})}[i][j] = 1$, of which the row index $i$ and the column index $j$ subject to $i \in [x_{1|n}, x_{2|n}]$ and $j \in [y_{1|n}, y_{2|n}]$.

Next, we obtain a masking map of the overall chest area, denoted $\mathcal{M}^{(\text{Chest})}$. We start by constructing the coordinates of the overall chest area by collecting the minimum values of $\{x_{1|n}\}_{n=1}^N$ and $\{y_{1|n}\}_{n=1}^N$, and the maximum values of $\{x_{2|n}\}_{n=1}^N$ and $\{y_{2|n}\}_{n=1}^N$, to construct $(\text{x}_1^{(\text{Chest})}, \text{x}_2^{(\text{Chest})}, \text{y}_1^{(\text{Chest})}, \text{y}_2^{(\text{Chest})})$. Similarly to $\mathcal{M}_n^{(\text{Ana})}$, an all-zero matrix is created, denoted $\mathcal{M}^{(\text{Chest})} \in \mathbb{R}^{\mathcal{W} \times \mathcal{H}}$, where $\mathcal{M}^{(\text{Chest})}[i][j] = 1$ if $i$ and $j$ subject to $i \in [\text{x}_1^{(\text{Chest})}, \text{x}_2^{(\text{Chest})}]$ and $j \in [\text{y}_1^{(\text{Chest})}, \text{y}_2^{(\text{Chest})}]$. The pseudo-code of the corresponding algorithm is shown in Alg. 1.

## A.2. Constructing Hierarchical Diagnostic Descriptions

To obtain three-level hierarchical diagnostic descriptions $\{H_m\}_{m=0}^{M=2}$, at each level we construct the corresponding labels of the diagnostic items and then construct the diagnostic descriptions using a few curated templates with blanks filled with the constructed annotations.

To begin with, we make use of abnormality annotations provided by Chest ImaGemone as the core abnormality labels to construct the labels of diagnostic items at each level. To obtain abnormality labels $\{l_{i|(n,k)}^{(\text{Abn})}\}$, for each $n^{th}$ anatomical part involved in the $i^{th}$ report, we collect the associated annotations of the pre-defined $K$ abnormalities, and then construct the labels of abnormality presentation, denoted as $\{l_{i|(n,k)}^{(\text{Abn})}\}_{k=1}^K$. $l_{i|(n,k)}^{(\text{Abn})} = 1$ indicates that the $k^{th}$ abnormality is presented on the $n^{th}$ anatomical parts; and $l_{i|(n,k)}^{(\text{Abn})} = 0$ otherwise.

Consequently, we construct the labels of diagnostic items at each level. The implementation details of the three-level diagnostic descriptions are presented as follows:

i) $H_0 \in \mathbb{R}^{N \times L}$ contains $N$ diagnostic items which describe whether the size and location of the $N$ anatomical parts are normal or not.
*Implementation*: We construct the labels of location normality $\{\{l_{i|n}^{(\text{loc})}\}_{n=1}^N\}$ and size normality $\{\{l_{i|n}^{(\text{size})}\}_{n=1}^N\}$ of the $N$ anatomical parts by their bounding box coordinates of $\{\{(\text{x}_{i|n}, \text{y}_{i|n}, \text{w}_{i|n}, \text{h}_{i|n})\}_{n=1}^N\}$ as described in Alg. 2.
ii) $H_1 \in \mathbb{R}^{N \times L}$ contains $N$ diagnostic items which describe whether the $N$ anatomical part are normal and with some medical devices or not;
*Implementation*: For each $i^{th}$ report, we construct the labels of anatomical normality $\{\{l_{i|n}^{(\text{nor})}\}_{n=1}^N\}$ and medical device $\{\{l_{i|n}^{(\text{dev})}\}_{n=1}^N\}$ of the $N$ anatomical parts.

To obtain the labels of anatomical normality $l_{i|n}^{(\text{nor})}$ of $n^{th}$ anatomical part, we make use of $\{l_{i|(n,k)}^{(\text{Abn})}\}_{k=1}^K$ to construct

$$l_{i|n}^{(\text{nor})} = \begin{cases} 0 & \exists k \in [1, K], l_{i|(n,k)}^{(\text{Abn})} = 1; \\ 1 & otherwise. \end{cases} \quad (1)$$

We also obtain $l_{i|n}^{(\text{dev})}$ by setting $l_{i|n}^{(\text{dev})} = 1$ if there exist the associated sentences of $n^{th}$ anatomical part which describe the placement of any medical device, tubes or lines; and $l_{i|n}^{(\text{dev})} = 0$ otherwise.
iii) $H_2 \in \mathbb{R}^{NK \times L}$ contains $NK$ diagnostic items indicating if a pre-defined set of $K$ abnormalities are detected in $N$ anatomical parts or not.
*Implementation*: We use $l_{i|(n,k)}^{(\text{abn})}$ directly as the labels of the diagnostic items at the level $m = 2$ for each report.

**Algorithm 1:** Construct masking maps of anatomical parts and chest area

---

**Input:** The bounding box coordinate sets $\{(x_n, y_n, w_n, h_n)\}_{n=1}^N$, the height and width of the image $(H, W)$, the height and width of visual feature map $(\mathcal{H}, \mathcal{W})$

**Output:** The masking map of anatomical parts $\{\mathcal{M}_n^{(\text{Ana})}\}_{n=1}^N$ and chest area $\mathcal{M}^{\text{Chest}}$

1 **def** *ReScale((x, y, w, h), (H, W), (\mathcal{H}, \mathcal{W}))***:**
2     $x_1 = x * \mathcal{W}/W$;
3     $x_2 = (x + w) * \mathcal{W}/W$;
4     $y_1 = y * \mathcal{H}/H$;
5     $y_2 = (y + h) * \mathcal{H}/H$;
6     **return** $(x_1, x_2, y_1, y_2)$;
7 **for** *each $n \in N$* **do**
8     Initial all-zero matrix $\mathcal{M}_n^{(\text{Ana})} \in \mathbb{R}^{\mathcal{W} \times \mathcal{H}}$;
9     $(x_{1|n}, x_{2|n}, y_{1|n}, y_{2|n}) \leftarrow$ ReScale$((x_n, y_n, w_n, h_n), (H, W), (\mathcal{H}, \mathcal{W}))$;
10     **for** *each $i \in [x_{1|n}, x_{2|n}]$* **do**
11        **for** *each $j \in [y_{1|n}, y_{2|n}]$* **do**
12           $\mathcal{M}_n^{(\text{Ana})}[i][j] = 1$
13 $x^{(\text{Chest})}, y^{(\text{Chest})} = \min\left(\{x_n\}\right)_{n=1}^N, \min\left(\{y_n\}\right)_{n=1}^N$;
14 $w^{(\text{Chest})}, h^{(\text{Chest})} = \max\left(\{x_n + w_n - x^{(\text{Chest})}\}\right)_{n=1}^N, \max\left(\{y_n + h_n - h^{(\text{Chest})}\}\right)_{n=1}^N$;
15 Initial all-zero matrix $\mathcal{M}^{(\text{Chest})} \in \mathbb{R}^{\mathcal{W} \times \mathcal{H}}$;
16 $(x_1^{(\text{Chest})}, x_2^{(\text{Chest})}, y_1^{(\text{Chest})}, y_2^{(\text{Chest})}) =$ ReScale$((x^{(\text{Chest})}, y^{(\text{Chest})}, w^{(\text{Chest})}, h^{(\text{Chest})}), (\mathcal{H}, \mathcal{W}))$;
17 **for** *each $i \in [x_1^{(\text{Chest})}, x_2^{(\text{Chest})}]$* **do**
18     **for** *each $j \in [y_1^{(\text{Chest})}, y_2^{(\text{Chest})}]$* **do**
19        $\mathcal{M}^{(\text{Chest})}[i][j] = 1$
20 **return** $\{\mathcal{M}_n^{(\text{Ana})}\}_{n=1}^N, \mathcal{M}^{(\text{Chest})}$;

---

Hierarchical diagnostic descriptions $\{H_m\}$ are then constructed by the labels of diagnostic items and a set of predefined templates, as illustrated in Table 1. We note that the diversity of the curated templates could affect the performance of the proposed model to be learned. However, using manual-design templates is one available way to construct hierarchical diagnostic descriptions. The intelligent and optimal construction of diagnostic descriptions would be systematically studied in future research.

## B. Model Implementation

The proposed model is implemented by PyTorch [6]. with a Python Package *Transformer* (`https://huggingface.co/`) adopted for the CLIP model. It takes 20 hours and 31 hours to fine-tune the `CLIP(SapBERT)` and learn our proposed model `AHIVE` on one NVIDIA 80G A100 GPU, respectively. For inference, it takes on average 0.398s to retrieve one report from 1000 candidates using the same GPU server.

All baselines are implemented according to their official codes and pre-trained weights[1]. The source code will be

---

[1]CLIP [7]:`https://huggingface.co/openai/`

released after the review process.

**Evaluation metric** The abnormality classes evaluated in the evaluation metric are detailed as follows:

`CE(11)` covers the top-11 most frequent abnormalities: *Lung opacity*, *Mass/Nodule*, *Airspace opacity*, *Pleural Effusion*, *Atelectasis*, Pneumothorax, *Medical device/tubelines*, *Edema*, *Consolidation*, *Enlarged cardiac silhouette*, *Lung lesion*.

`CE(11/5)` covers the top-11 abnormalities for each of the following 5 anatomical parts: *Left lung*, *Right lung*, *Left hilar structures*, *Right hilar structures* and *Cardiac silhouette*.

`CE(13+NL)` covers 13 abnormalities considered in [3]: *Enlarged cardiomediastinum*, *Cardiomegaly*, *Lesion*, *Lung opacity*, *Edema*, *Consolidation*, *Pneumonia*, *Atelectasis*, *Pneumothorax*, *Pleural effusion*, *Other*, *Fracture*, *Support*

---

`clip-vit-large-patch14`, SapBERT [5]: `https://huggingface.co/cambridgeltl/SapBERT-from-PubMedBERT-fulltext`, CXR-RePaiR [2]: `https://github.com/rajpurkarlab/CXR-RePaiR`, Mde-CLIP [8]:`https://github.com/RyanWangZf/MedCLIP`, BiomedCLIP [11]:`https://huggingface.co/microsoft/BiomedCLIP-PubMedBERT_256-vit_base_patch16_224`, X-REM [4]:`https://github.com/rajpurkarlab/X-REM`

| Template | $H_0$: "(1) is (2) the normal location. (1) is (3) the normal size."<br>$H_1$: "(1) is (4) abnormal. There is (5) medical device in (1)."<br>$H_2$: "There is (6) (7) in (1)." |
|---|---|
| Terms to fill | (1)=the term of $n^{th}$ anatomical part<br>(2)="in" if $l_n^{(\text{loc})} = 0$ otherwise "out of"<br>(3)="in" if $l_n^{(\text{size})} = 0$ otherwise "out of"<br>(4)="in" if $l_n^{(\text{nor})} = 0$ otherwise "not in"<br>(5)="no" if $l_n^{(\text{dev})} = 0$ otherwise keep blank<br>(6)="no" if $l_{n,k}^{(\text{abn})} = 0$ otherwise keep blank<br>(7)=the term of $k^{th}$ abnormality on $n^{th}$ anatomical part |

Table 1. An example of templates with blanks for constructing the hierarchical diagnostic descriptions, where (#ID) represents the blank to be filled according to the labels of the diagnostic items.

*devices*, and NL stands for *Normality* [1].

# References

[1] Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. Generating radiology reports via memory-driven transformer. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1439–1449, 2020. 3

[2] Mark Endo, Rayan Krishnan, Viswesh Krishna, Andrew Y Ng, and Pranav Rajpurkar. Retrieval-based chest x-ray report generation using a pre-trained contrastive language-image model. In *Machine Learning for Health*, pages 209–219. PMLR, 2021. 2

[3] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the Conference of Association for the Advance of Artificial Intelligence*, pages 590–597, 2019. 2

[4] Jaehwan Jeong, Katherine Tian, Andrew Li, Sina Hartung, Subathra Adithan, Fardad Behzadi, Juan Calle, David Osayande, Michael Pohlen, and Pranav Rajpurkar. Multimodal image-text matching improves retrieval-based chest x-ray report generation. In *Medical Imaging with Deep Learning*, 2023. 2

[5] Fangyu Liu, Ehsan Shareghi, Zaiqiao Meng, Marco Basaldella, and Nigel Collier. Self-alignment pretraining for biomedical entity representations. *arXiv preprint arXiv:2010.11784*, 2020. 2

[6] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing systems*, 32, 2019. 2

[7] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 2

[8] Zifeng Wang, Zhenbang Wu, Dinesh Agarwal, and Jimeng Sun. Medclip: Contrastive learning from unpaired medical images and text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 3876–3887, 2022. 2

[9] Joy T Wu, Nkechinyere N Agu, Ismini Lourentzou, Arjun Sharma, Joseph A Paguio, Jasper S Yao, Edward C Dee, William Mitchell, Satyananda Kashyap, Andrea Giovannini, et al. Chest imagenome dataset for clinical reasoning. *arXiv preprint arXiv:2108.00316*, 2021. 1

[10] Joy T Wu, Nkechinyere N Agu, Ismini Lourentzou, Arjun Sharma, Joseph A Paguio, Jasper S Yao, Edward C Dee, William Mitchell, Satyananda Kashyap, Andrea Giovannini, et al. Chest imagenome dataset (version 1.0. 0). *PhysioNet*, 5:18, 2021. 1

[11] Sheng Zhang, Yanbo Xu, Naoto Usuyama, Jaspreet Bagga, Robert Tinn, Sam Preston, Rajesh Rao, Mu Wei, Naveen Valluri, Cliff Wong, et al. Large-scale domain-specific pretraining for biomedical vision-language processing. *arXiv preprint arXiv:2303.00915*, 2023. 2

**Algorithm 2:** Construct labels of location and size normality

---

**Input:** The set of anatomical bounding box coordinate $\mathcal{B} = \{\{(\mathrm{x}_{i|n}, \mathrm{y}_{i|n}, \mathrm{w}_{i|n}, \mathrm{h}_{i|n})\}_{n=1}^{N}\}$, the set of normality labels $L_{(\text{normality})} = \{l_i^{(\text{nor})}\}$

**Output:** The sets of location normality labels $L_{(\text{location})} = \{\{l_{i|n}^{(\text{loc})}\}_{n=1}^{N}\}$ and size normality labels $L_{(\text{size})} = \{\{l_{i|n}^{(\text{size})}\}_{n=1}^{N}\}$

1 **def** *RelativeSize(*$\{(\mathrm{x}_n, \mathrm{y}_n, \mathrm{w}_n, \mathrm{h}_n)\}_{n=1}^{N}$*)***:**
2     $\mathrm{x}^{(\text{Chest})}, \mathrm{y}^{(\text{Chest})} = \min\left(\{\mathrm{x}_n\}\right)_{n=1}^{N}, \min\left(\{\mathrm{y}_n\}\right)_{n=1}^{N}$;
3     $\mathrm{w}^{(\text{Chest})}, \mathrm{h}^{(\text{Chest})} = \max\left(\{\mathrm{x}_n + \mathrm{w}_n - \mathrm{x}^{(\text{Chest})}\}\right)_{n=1}^{N}, \max\left(\{\mathrm{y}_n + \mathrm{h}_n - \mathrm{h}^{(\text{Chest})}\}\right)_{n=1}^{N}$;
4     **for** *each $n \in N$* **do**
5        $\tilde{s}_n = \mathrm{w}_n \cdot \mathrm{h}_n / (\mathrm{w}^{(\text{Chest})} \cdot \mathrm{h}^{(\text{Chest})})$;
6     **return** $(\mathrm{x}^{(\text{Chest})}, \mathrm{y}^{(\text{Chest})}, \mathrm{w}^{(\text{Chest})}, \mathrm{h}^{(\text{Chest})}), \{\tilde{s}_n\}_{n=1}^{N}$;

7 **def** *RelativeCenter(*$\{(\mathrm{x}_n, \mathrm{y}_n, \mathrm{w}_n, \mathrm{h}_n)\}_{n=1}^{N}, (\mathrm{x}^{(\text{Chest})}, \mathrm{y}^{(\text{Chest})}, \mathrm{w}^{(\text{Chest})}, \mathrm{h}^{(\text{Chest})})$*)***:**
8     **for** *each $n \in N$* **do**
9        $\tilde{x}_n = \mathrm{x}_n + \mathrm{w}_n/2 - (\mathrm{x}^{(\text{Chest})} + \mathrm{w}^{(\text{Chest})}/2)/\mathrm{w}^{(\text{Chest})}$;
10       $\tilde{y}_n = \mathrm{y}_n + \mathrm{h}_n/2 - (\mathrm{y}^{(\text{Chest})} + \mathrm{h}^{(\text{Chest})}/2)/\mathrm{h}^{(\text{Chest})}$;
11     **return** $\{(\tilde{x}_n, \tilde{y}_n)\}$;

12 **for** *each $n \in N$* **do**
13     Initial $X_n = \emptyset, Y_n = \emptyset, S_n = \emptyset$ ;

14 **for** *each $\{(\mathrm{x}_{i|n}, \mathrm{y}_{i|n}, \mathrm{w}_{i|n}, \mathrm{h}_{i|n})\}_{n=1}^{N} \in \mathcal{B}$* **do**
15     **if** $l_i^{(\text{nor})} = 1$ **then**
16        $(\mathrm{x}_i^{(\text{Chest})}, \mathrm{y}_i^{(\text{Chest})}, \mathrm{w}_i^{(\text{Chest})}, \mathrm{h}_i^{(\text{Chest})}), \{\tilde{s}_{i|n}\}_{n=1}^{N} \leftarrow \text{RelativeSize}(\{(\mathrm{x}_{i|n}, \mathrm{y}_{i|n}, \mathrm{w}_{i|n}, \mathrm{h}_{i|n})\}_{n=1}^{N})$;
17        $\{(\tilde{x}_{i|n}, \tilde{y}_{i|n})\}_{n=1}^{N} \leftarrow \text{RelativeCenter}(\mathrm{x}_i^{(\text{Chest})}, \mathrm{y}_i^{(\text{Chest})}, \mathrm{w}_i^{(\text{Chest})}, \mathrm{h}_i^{(\text{Chest})}), \{(\mathrm{x}_{i|n}, \mathrm{y}_{i|n}, \mathrm{w}_{i|n}, \mathrm{h}_{i|n})\}_{n=1}^{N})$;
18        **for** *each $n \in N$* **do**
19           Add $\tilde{x}_{i|n}$ to $X_n$, add $\tilde{y}_{i|n}$ to $Y_n$, add $\tilde{s}_{i|n}$ to $S_n$;

20 **for** *each $n \in N$* **do**
21     $\mu_{x|n}, \sigma_{x|n} \leftarrow \text{mean-and-std}(X_n)$;
22     $\mu_{y|n}, \sigma_{y|n} \leftarrow \text{mean-and-std}(Y_n)$;
23     $\mu_{s|n}, \sigma_{s|n} \leftarrow \text{mean-and-std}(S_n)$;

24 Initial $L_{(\text{location})} = \emptyset, L_{(\text{location})} = \emptyset$;
25 **for** *each $\{(\mathrm{x}_{i|n}, \mathrm{y}_{i|n}, \mathrm{w}_{i|n}, \mathrm{h}_{i|n})\}_{n=1}^{N} \in \mathcal{B}$* **do**
26     $(\mathrm{x}_i^{(\text{Chest})}, \mathrm{y}_i^{(\text{Chest})}, \mathrm{w}_i^{(\text{Chest})}, \mathrm{h}_i^{(\text{Chest})}), \{\tilde{s}_{i|n}\}_{n=1}^{N} \leftarrow \text{RelativeSize}(\{(\mathrm{x}_{i|n}, \mathrm{y}_{i|n}, \mathrm{w}_{i|n}, \mathrm{h}_{i|n})\}_{n=1}^{N})$;
27     $\{(\tilde{x}_{i|n}, \tilde{y}_{i|n})\}_{n=1}^{N} \leftarrow \text{RelativeCenter}(\mathrm{x}_i^{(\text{Chest})}, \mathrm{y}_i^{(\text{Chest})}, \mathrm{w}_i^{(\text{Chest})}, \mathrm{h}_i^{(\text{Chest})}), \{(\mathrm{x}_{i|n}, \mathrm{y}_{i|n}, \mathrm{w}_{i|n}, \mathrm{h}_{i|n})\}_{n=1}^{N})$;
28     **for** *each $n \in N$* **do**
29        Initialize $l_{i|n}^{(\text{loc})} = 0$ and $l_{i|n}^{(\text{size})} = 0$;
30        **if** $\tilde{x}_{i|n} \in (\mu_{x|n} - \sigma_{x|n}, \mu_{x|n} + \sigma_{x|n})$ *and* $\tilde{y}_n \in (\mu_{y|n} - \sigma_{y|n}, \mu_{y|n} + \sigma_{y|n})$ **then**
31           Set $l_{i|n}^{(\text{loc})} = 1$;
32        **if** $\tilde{s}_n \in (\mu_{s|n} - \sigma_{s|n}, \mu_{s|n} + \sigma_{s|n})$ **then**
33           Set $l_{i|n}^{(\text{size})} = 1$;
34        Add $l_{i|n}^{(\text{loc})}$ to $L_{(\text{location})}$, $l_{i|n}^{(\text{size})}$ to $L_{(\text{Size})}$;

35 **return** $L_{(\text{location})}, L_{(\text{size})}$;