

## – Supplementary Material –

# Foundation X: Integrating Classification, Localization, and Segmentation through Lock-Release Pretraining Strategy for Chest X-ray Analysis

## A. Experiment details

Here, we discuss the setup of the training process for Foundation X. The backbone is Swin-B, initialized with Ark-6 [22] pretrained weights. For the classification task, linear layers serve as classification heads. For localization, we integrate the DINO localization approach [38], modify to handle multiple datasets with one localization encoder and multiple localization decoders. For segmentation, we use UperNet [36], modify to include multiple segmentation heads. We pretrain Foundation X on all 11 datasets (see Table 1) using a single A100 GPU, employing the Cyclic and Lock-Release pretraining strategies. We also employ the Student-Teacher learning paradigm, where the teacher model is an exact copy of the student model at the start. The teacher model is updated after each epoch using an exponential moving average (EMA) [32] with a momentum of 0.80. The configuration for Foundation X is detailed in Table 7.

|                                      |                           |
|--------------------------------------|---------------------------|
| Backbone                             | Swin-B <sup>†</sup> [20]  |
| Classification Branch                | Linear Layer <sup>‡</sup> |
| Localization Branch                  | DINO <sup>‡</sup> [38]    |
| Segmentation Branch                  | UperNet <sup>‡</sup> [36] |
| Input Resolution                     | 224 x 224                 |
| Optimizer                            | AdamW                     |
| Batch Size                           | 24                        |
| Number of Workers                    | 12                        |
| Backbone Learning Rate               | 1e-5                      |
| Localization Learning Rate           | 1e-4                      |
| Segmentation Learning Rate           | 1e-4                      |
| Learning Rate Scheduler              | Step-decay                |
| Evaluation Metric for Classification | AUC                       |
| Evaluation Metric for Localization   | mAP40                     |
| Evaluation Metric for Segmentation   | Dice                      |

<sup>†</sup> Initialized with Ark-6 [22] pretrained weights.

<sup>‡</sup> Initialized with random weights.

Table 7. Experiment settings for Foundation X.

## B. Model Parameters

The Foundation X model consists of several key components, contributing to a total of approximately 173 million parameters (See Table 9). The backbone, responsible for feature extraction, accounts for 86.8 million parameters. The localization encoder adds 7.7 million parameters, while the localization decoders total 57.6 million, with each decoder contributing 9.6 million. The segmentation decoder

comprises 20.9 million parameters. Although adding a dedicated localization decoder for each dataset increases the model size, only the relevant decoder is active during training, with the others, along with the classification and segmentation heads, remaining frozen. This approach keeps the computational load manageable and ensures efficient GPU utilization.

## C. Lock-Release pretraining strategy

Foundation X effectively handles classification, localization, and segmentation tasks. The model leverages the Student-Teacher learning paradigm along with Cyclic and Lock-Release pretraining strategies, ensuring it retains general knowledge for all tasks while avoiding overfitting to any single task. In Table 8, we illustrate the Lock-Release pretraining strategy using the VinDr-CXR organ dataset for organ localization and segmentation. For this demonstration, we treat localization and segmentation of the heart, left lung, and right lung as separate tasks.

## D. Cross-Dataset and Cross-Task learning analysis

The full figure illustrating the Cross-Dataset and Cross-Task learning analysis for all six datasets is included in this supplementary material (see Figure 4). This figure highlights the performance trends of Foundation X across various datasets under both focused and unfocused training scenarios, showcasing its ability to generalize and retain knowledge effectively through the Cyclic and Lock-Release pretraining strategies. We include plots for these six datasets because they contain multiple tasks, including classification, localization, and segmentation.

## E. Ablation study

The ablation studies demonstrate the effectiveness of the Student-Teacher learning paradigm, Cyclic and Lock-Release pretraining strategies across various tasks. Foundation X-L (see Table 11), trained on six disease localization tasks, generally outperforms the baseline model Swin-B + DINO. Similarly, Foundation X-S (see Table 12), trained on three disease segmentation datasets, consistently surpasses the baseline model Swin-B + UperNet. Additionally, Foundation X-CL (see Table 13), which handles both classification and localization tasks, and Foundation X-LS (see Table 14), which integrates localization and segmentation tasks, both show superior performance compared to their

---

**Algorithm 1:** A round of Foundation X’s Cyclic Lock-Release pretraining

---

**Data:** Datasets:  $\mathcal{D} = \{D_1, D_2, \dots, D_n\}$ ; Sample: image-label pair  $(x, y) \in \mathcal{D}_i$   
**Functions:** Data augmentation:  $\varepsilon$ ; Dataset/task-specific losses:  $\{\mathcal{LD}1(\cdot, \cdot), \mathcal{LD}2(\cdot, \cdot), \dots, \mathcal{LD}n(\cdot, \cdot)\}$ ; Consistency loss:  $\{\mathcal{Lconst}(\cdot, \cdot)\}$ ; Loss update by AdamW optimizer:  $Update_{adamw}(\cdot, \cdot)$   
**Trainable Parameters:** Student’s encoder, localization encoder, segmentation decoder:  $e_s, LocEnc_s, SegDec_s$ ; Classification heads  $\mathcal{C} = \{C_1, C_2, \dots, C_n\}$ ; Localization decoders  $\mathcal{L} = \{L_1, L_2, \dots, L_n\}$ ; Segmentation heads  $\mathcal{S} = \{S_1, S_2, \dots, S_n\}$ ;  
**Stop Gradient:** Teacher’s encoder, localization encoder, segmentation decoder:  $e_t, LocEnc_t, SegDec_t$ ;  
**Hyperparameters:** Momentum:  $\lambda$

```
1  $\{e_t, LocEnc_t, SegDec_t\} \leftarrow \{e_s, LocEnc_s, SegDec_s\}$  // initialize teacher with student’s parameters
2 for  $D_i$  in  $D_1, D_2, \dots, D_n$  do
  /* train student for one epoch */
3   for  $(x, y)$  in  $D_i$  do
4      $x' = \varepsilon(x)$ 
5     if  $D_i$  has Classification Annotation then
6       for  $j \leftarrow 1$  to 2 do
7         if  $j = 1$  then
8           Freeze  $\{e_s\}$  // Lock mode on, using a random half of the dataset
9         else
10          Unfreeze  $\{e_s\}$  // Release mode on, using full dataset
11           $emb_t, emb_s = e_t(x'), e_s(x')$ 
12           $pred = C_i(emb_s)$ 
13           $Loss = \mathcal{LD}i(pred, y) + \mathcal{Lconst}1(emb_t, emb_s)$ 
14          Update( $\{e_s, p_s, C_i\}, Loss$ )
15     if  $D_i$  has Localization Annotation then
16       for  $j \leftarrow 1$  to 2 do
17         if  $j = 1$  then
18           Freeze  $\{e_s, LocEnc_s\}$  // Lock mode on, using a random half of the dataset
19         else
20           Unfreeze  $\{e_s, LocEnc_s\}$  // Release mode on, using full dataset
21            $emb_t, emb_s = e_t(x'), e_s(x')$ 
22            $embLocEnc_s, embLocEnc_t = LocEnc_s(emb_s), LocEnc_t(emb_t)$ 
23            $pred = L_i(embLocEnc_s)$ 
24            $Loss = \mathcal{LD}i(pred, y) + \mathcal{Lconst}(emb_t, emb_s) + \mathcal{Lconst}(embLocEnc_t, embLocEnc_s)$ 
25           Update( $\{e_s, LocEnc_s, L_i\}, Loss$ )
26     if  $D_i$  has Segmentation Annotation then
27       for  $j \leftarrow 1$  to 2 do
28         if  $j = 1$  then
29           Freeze  $\{e_s, SegDec_s\}$  // Lock mode on, using a random half of the dataset
30         else
31           Unfreeze  $\{e_s, SegDec_s\}$  // Release mode on, using full dataset
32            $emb_t, emb_s = e_t(x'), e_s(x')$ 
33            $embSegDec_s, embSegDec_t = SegDec_s(emb_s), SegDec_t(emb_t)$ 
34            $pred = S_i(embSegDec_s)$ 
35            $Loss = \mathcal{LD}i(pred, y) + \mathcal{Lconst}(emb_t, emb_s) + \mathcal{Lconst}(embSegDec_t, embSegDec_s)$ 
36           Update( $\{e_s, SegDec_s, S_i\}, Loss$ )
37   /* Update teacher by student’s parameters via epoch-wise EMA */
    $\{e_t, LocEnc_t, SegDec_t\} \leftarrow \lambda\{e_t, LocEnc_t, SegDec_t\} + (1 - \lambda)\{e_s, LocEnc_s, SegDec_s\}$ 
```

---

|         | Epoch # | Data Size | Backbone | Loc.Enc | Loc.Dec | Seg.Dec | Seg.Head | Mode    | Training Task              |
|---------|---------|-----------|----------|---------|---------|---------|----------|---------|----------------------------|
| Cycle 1 | 1       | Half      | F        | F       | T       | -       | -        | Lock    | Localization of Heart      |
|         | 2       | Full      | T        | T       | T       | -       | -        | Release | Localization of Heart      |
|         | 3       | Half      | F        | F       | T       | -       | -        | Lock    | Localization of Left Lung  |
|         | 4       | Full      | T        | T       | T       | -       | -        | Release | Localization of Left Lung  |
|         | 5       | Half      | F        | F       | T       | -       | -        | Lock    | Localization of Right Lung |
|         | 6       | Full      | T        | T       | T       | -       | -        | Release | Localization of Right Lung |
|         | 7       | Half      | F        | -       | -       | F       | T        | Lock    | Segmentation of Heart      |
|         | 8       | Full      | T        | -       | -       | T       | T        | Release | Segmentation of Heart      |
|         | 9       | Half      | F        | -       | -       | F       | T        | Lock    | Segmentation of Left Lung  |
|         | 10      | Full      | T        | -       | -       | T       | T        | Release | Segmentation of Left Lung  |
|         | 11      | Half      | F        | -       | -       | F       | T        | Lock    | Segmentation of Right Lung |
|         | 12      | Full      | T        | -       | -       | T       | T        | Release | Segmentation of Right Lung |

Table 8. Demonstrating the Lock-Release pretraining strategy for organ localization and segmentation using the VinDr-CXR dataset. The model completes a single cycle when it goes through all tasks once (from epoch #1 to #12). 'F' denotes a frozen component, and 'T' denotes a trainable component. In Lock mode, the model is trained using half of the dataset, while in Release mode, it is trained using the full dataset. After each epoch in Release mode, the model is tested on the localization and segmentation of the heart, left lung, and right lung.

| Component                        | Parameters         |
|----------------------------------|--------------------|
| Backbone                         | 86,751,673 [+]     |
| Classification Heads             | 70,725 [+]         |
| Localization Encoder             | 7,693,056 [+]      |
| Localization Decoders            | 57,653,868 [+]     |
| <i>Each Localization Decoder</i> | 9,608,978          |
| Segmentation Decoder             | 20,894,464 [+]     |
| Segmentation Heads               | 20,754 [+]         |
| <b>Total</b>                     | <b>173,084,540</b> |

Table 9. Parameter distribution across the key components of the Foundation X model, trained on 11 datasets and 20 tasks.

respective baseline methods in most cases.

The ablation study presented in Table 10 highlights the impact of incorporating the Lock-Release pretraining strategy and the Student-Teacher learning paradigm on the performance of the Foundation X model. The results demonstrate that when both components are enabled, the model achieves the best performance across all evaluated VinDr-CXR organs (Heart, Left Lung, and Right Lung) localization. Specifically, the combination of Lock-Release and Student-Teacher results in the highest mAP, with scores of 88.39% for Heart, 95.78% for Left Lung, and 96.78% for Right Lung. These findings suggest that each component complements the other, with the Lock-Release strategy preventing task-specific overfitting and the Student-Teacher paradigm ensuring stable learning by reducing drastic model shifts. Together, these strategies create a synergistic effect that enhances the model’s generalization and overall performance, outperforming configurations where one or both components are disabled. This highlights the importance of integrating both the Lock-Release strategy and the Student-Teacher paradigm to maximize the effectiveness of our approach.

| Lock-Release | Student-Teacher | Heart        | Left Lung    | Right Lung   |
|--------------|-----------------|--------------|--------------|--------------|
| ✗            | ✗               | 85.45        | 93.63        | 94.47        |
| ✗            | ✓               | 86.41        | 94.39        | 94.95        |
| ✓            | ✗               | 87.50        | 95.37        | 96.44        |
| ✓            | ✓               | <b>88.39</b> | <b>95.78</b> | <b>96.78</b> |

Table 10. Ablation study is conducted on the VinDr-CXR organ localization dataset. We evaluate the model with and without the Lock-Release pretraining strategy, as well as with and without the Student-Teacher model. The results demonstrate that the Foundation X model achieves comparatively better performance when both the Lock-Release pretraining strategy and the Student-Teacher learning paradigm are employed.

## F. Acknowledgements

This research was partially supported by ASU and Mayo Clinic through a Seed Grant and an Innovation Grant, as well as by the NIH under Award Number R01HL128785. The authors are solely responsible for the content, which does not necessarily reflect the official views of the NIH. This work also utilized GPUs provided by ASU Research Computing, Bridges-2 at the Pittsburgh Supercomputing Center (allocated under BCS190015), and Anvil at Purdue University (allocated under MED220025). These resources are supported by the Advanced Cyberinfrastructure Coordination Ecosystem: Services & Support (ACCESS) program, funded by the National Science Foundation under grants #2138259, #2138286, #2138307, #2137603, and #2138296. We also extend our gratitude to Anirudh Kaniyar Narayana Iyengar for his contributions to collecting localization data, preparing bounding boxes in COCO format, and developing some of the dataloaders. Finally, the content of this paper is covered by patents pending.

| Dataset        | Baseline Loc. <sup>†</sup><br>[mAP40%] | Foundation X-L<br>[mAP40%] |
|----------------|--|----------------------------|
| TBX11K         | <b>78.10</b>                           | 77.77 ↓ 0.33               |
| NODE21         | 37.50                                  | <b>42.79</b> ↑ 5.29        |
| CANDID-PTX     | 50.90                                  | <b>53.75</b> ↑ 2.85        |
| RSNA Pneumonia | 21.70                                  | <b>29.37</b> ↑ 7.67        |
| ChestX-Det     | 38.00                                  | <b>40.13</b> ↑ 2.13        |
| SIIM-ACR       | 28.00                                  | <b>36.20</b> ↑ 8.20        |

<sup>†</sup> Swin-B version of DINO where the backbone is initialized with Ark-6 pretrained weights.

Table 11. We train Foundation X-L on six disease localization tasks utilizing Cyclic and Lock-Release pretraining strategies and compare its performance with the baseline model, DINO [38]. In most cases, Foundation X-L outperforms the baseline across the datasets during pretraining.

| Dataset    | Baseline Seg. <sup>†</sup><br>[Dice%] | Foundation X-S<br>[Dice%] |
|------------|---------------------------------------|---------------------------|
| CANDID-PTX | 86.36                                 | <b>89.58</b> ↑ 3.23       |
| ChestX-Det | 79.33                                 | <b>83.46</b> ↑ 4.13       |
| SIIM-ACR   | 81.92                                 | <b>83.83</b> ↑ 1.91       |

<sup>†</sup> Swin-B version of UperNet where the backbone is initialized with Ark-6 pretrained weights.

Table 12. We train Foundation X-S on three disease segmentation datasets using the Cyclic and Lock-Release pretraining strategies and compare its performance with the baseline model, UperNet [36]. In all cases, Foundation X-S outperforms the baseline across the datasets during pretraining.

| Dataset        | Baseline Cls.      | Baseline Loc. | Foundation X-CL     |                      |
|----------------|--------------------|---------------|---------------------|----------------------|
|                | [AUC%]             | [mAP40%]      | [AUC%]              | [mAP40%]             |
| TBX11K         | 99.89±0.06         | <b>78.10</b>  | <b>99.96</b> ↑ 0.07 | 72.56 ↓ 5.54         |
| NODE21         | 99.35±0.45         | 37.50         | <b>99.68</b> ↑ 0.33 | <b>47.54</b> ↑ 10.04 |
| CANDID-PTX     | 72.61±0.57         | 50.90         | <b>74.00</b> ↑ 1.39 | <b>51.61</b> ↑ 0.71  |
| RSNA Pneumonia | 88.87±0.21         | 21.70         | <b>96.57</b> ↑ 7.70 | <b>26.08</b> ↑ 4.38  |
| ChestX-Det     | <b>88.17</b> ±0.33 | <b>38.00</b>  | 81.82 ↓ 6.35        | 37.03 ↓ 0.97         |
| SIIM-ACR       | 95.01±0.16         | 28.00         | <b>95.19</b> ↑ 0.18 | <b>34.98</b> ↑ 6.98  |

Table 13. We train Foundation X-CL on six disease datasets, each containing both classification and localization annotations, using our Cyclic and Lock-Release pretraining strategies. The table demonstrates that, in most cases, Foundation X-CL outperforms the baseline methods during pretraining.

| Dataset        | Baseline Loc. | Baseline Seg. | Foundation X-LS     |                     |
|----------------|---------------|---------------|---------------------|---------------------|
|                | [mAP40%]      | [Dice%]       | [mAP40%]            | [Dice%]             |
| TBX11K         | <b>78.10</b>  | -             | 73.03 ↓ 5.07        | -                   |
| NODE21         | 37.50         | -             | <b>46.09</b> ↑ 8.59 | -                   |
| CANDID-PTX     | 50.90         | 86.36         | <b>53.01</b> ↑ 2.11 | <b>89.47</b> ↑ 3.11 |
| RSNA Pneumonia | 21.70         | -             | <b>27.80</b> ↑ 6.10 | -                   |
| ChestX-Det     | 38.00         | <b>79.33</b>  | <b>39.22</b> ↑ 1.22 | 70.90 ↓ 8.43        |
| SIIM-ACR       | 28.00         | 81.92         | <b>36.63</b> ↑ 8.63 | <b>84.25</b> ↑ 2.33 |

Table 14. We train Foundation X-LS on six disease localization and three disease segmentation datasets, using our Cyclic and Lock-Release pretraining strategies. The table demonstrates that, in most cases, Foundation X-LS outperforms the baseline methods during pretraining.

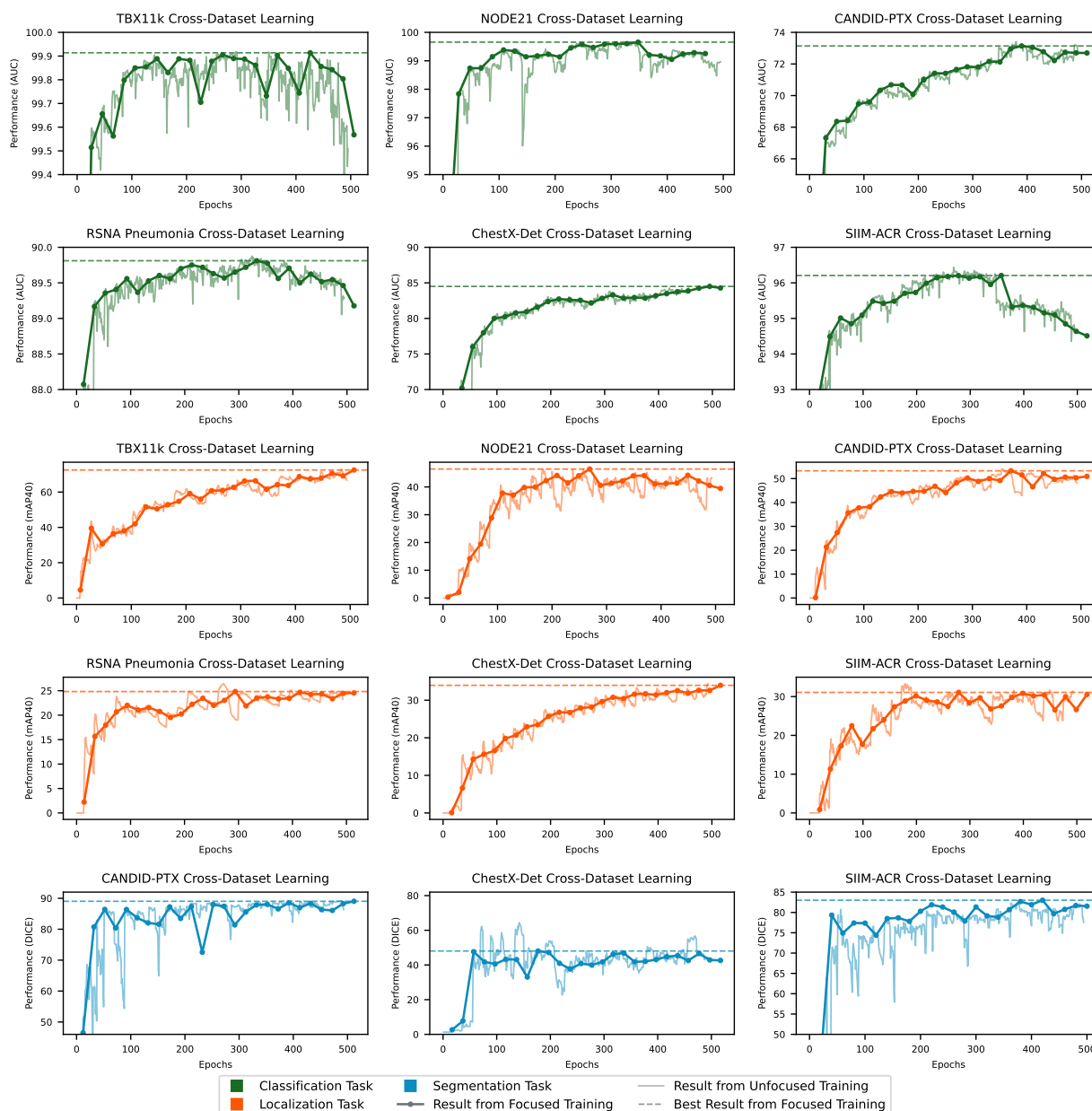


Figure 4. Cross-Dataset & Cross-Task Learning Analysis. The figure demonstrates the performance trends of Foundation X across multiple datasets for both focused and unfocused training scenarios. Focused training refers to scenarios where the model is explicitly trained on the specific dataset being evaluated, while unfocused training refers to scenarios where the model is trained on other datasets and not directly on the dataset being evaluated. The green, orange, and blue lines represent the classification, localization, and segmentation tasks, respectively. Dark-colored lines indicate the testing results during focused training, where the model is explicitly trained on the specific dataset. Light-colored lines show the testing results during unfocused training, where the model is trained on other datasets but tested on the specific dataset. Dashed lines represent the best testing results achieved from focused training for each specific dataset. The results indicate that, during unfocused training, initial performance dips are common as the model is not explicitly trained on the specific dataset. However, performance improves over time, demonstrating the model’s ability to generalize effectively and retain knowledge due to the Cyclic and Lock-Release pretraining strategies. In all cases, the unfocused training results do not drift away from the task, highlighting the model’s efficient generalization and knowledge retention. Additionally, in some instances, unfocused training achieves even better performance than focused training, showcasing the advantages of cross-task and cross-dataset learning in enhancing the overall capabilities of Foundation X.

| Task          | Dataset                | Official Split | Train Split             | Val Split     | Test Split   | Expert Labels   |  |
|---------------|------------------------|----------------|-------------------------|---------------|--------------|---|--|
| CLS           | CheXpert [9]           | ✓              | 223415                  | 234           | -            | No finding, Enlarged Cardiomeastinum, Cardiomegaly, Lung Opacity, Lung Lesion, Edema, Consolidation, Pneumonia, Atelectasis, Pneumothorax, Pleural Effusion, Pleural Other, Fracture, Support Devices           |  |
|               | NIH ChestX-ray14 [35]  | ✓              | 75312                   | 11212         | 25596        | Atelectasis, Cardiomegaly, Effusion, Infiltration, Mass, Nodule, Pneumonia, Pneumothorax, Consolidation, Edema, Emphysema, Fibrosis, Pleural Thickening, Hernia   |  |
|               | VinDr-CXR [25]         | ✓              | 15000                   | -             | 3000         | PE, Lung Tumor, Pneumonia, Tuberculosis, Other diseases, No finding   |  |
|               | NIH Shenzhen CXR [11]  | ✗              | 463                     | 65            | 134          | Tuberculosis  |  |
|               | MIMIC-II [12]          | ✓              | 368878                  | 2991          | 5159         | No finding, Enlarged Cardiomeastinum, Cardiomegaly, Lung Opacity, Lung Lesion, Edema, Consolidation, Pneumonia, Atelectasis, Pneumothorax, Pleural Effusion, Pleural Other, Fracture, Support Devices           |  |
| CLS, LOC      | TBX11k [19]            | ✓              | 6600                    | 1800          | -            | Tuberculosis  |  |
|               | NODE21 [31]            | ✗              | 4178                    | -             | 1046         | Nodule  |  |
|               | RSNA Pneumonia [26]    | ✓              | 21295                   | 2680          | 2709         | CLS: No lung opacity/Not normal, Normal, Lung Opacity; LOC: Pneumonia   |  |
| CLS, LOC, SEG | CANDID-PTX [3]         | ✓              | 13748                   | 1964          | 3928         | Pneumothorax  |  |
|               | ChestX-Det [16]        | ✓              | 3025                    | -             | 553          | Atelectasis, Calcification, Cardiomegaly, Consolidation, Diffuse Nodule, Effusion, Emphysema, Fibrosis, Fracture, Mass, Nodule, Pleural Thickening, Pneumothorax  |  |
|               | SIIM-ACR [1]           | ✗              | 9607                    | 1068          | 1372         | Pneumothorax  |  |
|               |                        |                | <b>Total CLS images</b> | <b>741521</b> | <b>22014</b> | <b>43997</b>  |  |
|               |                        |                | <b>Total LOC images</b> | <b>58453</b>  | <b>7512</b>  | <b>9608</b>   |  |
|               |                        |                | <b>Total SEG images</b> | <b>26380</b>  | <b>3032</b>  | <b>5853</b>   |  |
| LOC FT        | VinDr-CXR [25]         | ✓              | 15000                   | -             | 3000         | Aortic enlargement, Atelectasis, Calcification, Cardiomegaly, Consolidation, ILD, Infiltration, Lung Opacity, Nodule/Mass, Other lesion, Pleural effusion, Pleural thickening, Pneumothorax, Pulmonary fibrosis |  |
| SEG FT        | CheXmask VinDr-CXR [5] | ✓              | 15000                   | -             | 3000         | Heart, Left Lung, Right Lung  |  |
|               | VinDr-RibCXR [24]      | ✓              | 196                     | -             | 49           | 20 Ribs   |  |
|               | NIH Montgomery [11]    | ✗              | 92                      | 15            | 31           | Lung  |  |
|               | JSRT [33]              | ✗              | 173                     | 25            | 49           | Heart, Lung, Clavicle   |  |

Table 15. Foundation X was pretrained on the above 11 classification datasets, 6 localization datasets, and 3 segmentation datasets. CLS stands for classification task, LOC stands for localization task, SEG stands for segmentation task. "CLS, LOC" denotes the datasets used for classification and localization tasks. "CLS, LOC, SEG" denotes the datasets used for classification, localization, and segmentation tasks. "LOC FT" and "SEG FT" denotes the datasets used only during the finetuning of the localization and segmentation task, respectively.