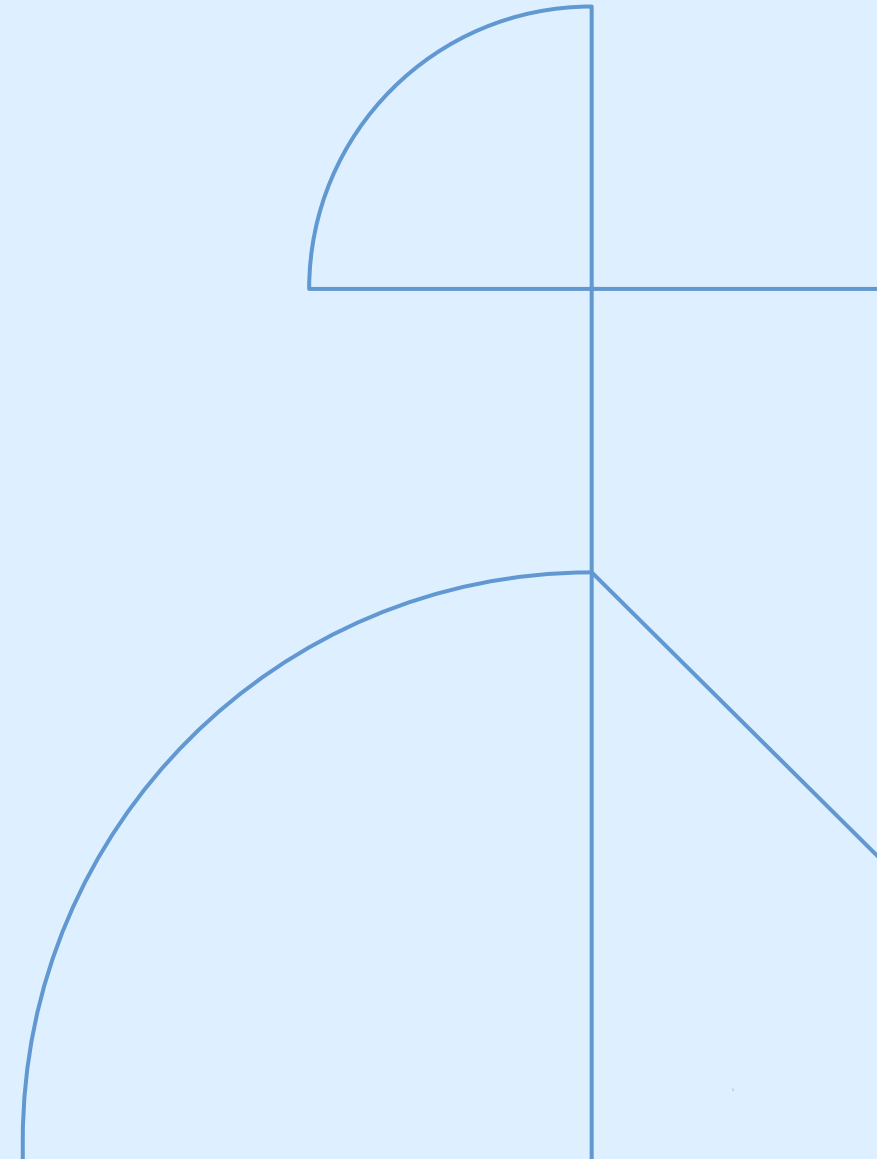# Learning from Offline Foundation Features with Tensor Augmentations

Emir Konuk, Christos Matsoukas, Moein Sorkhei, Phitchapha Lertsiravarameth, Kevin Smith
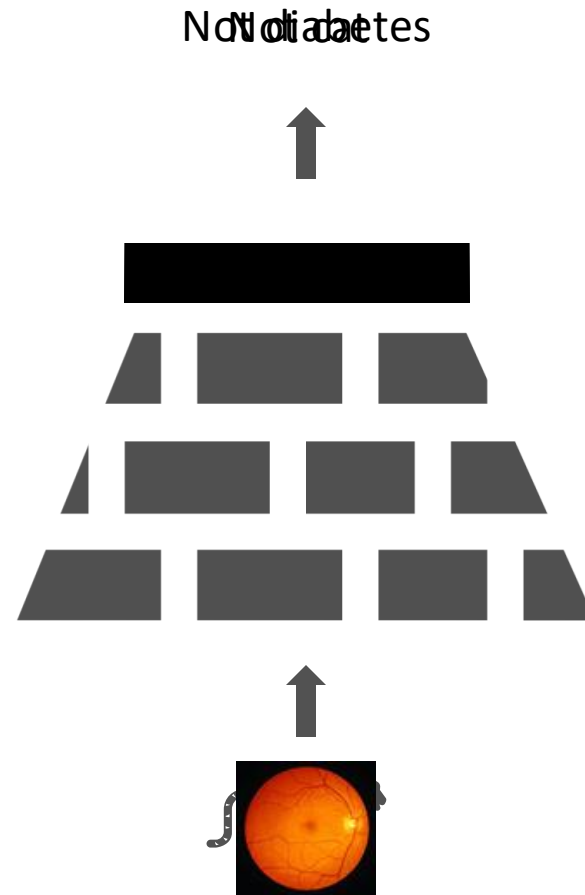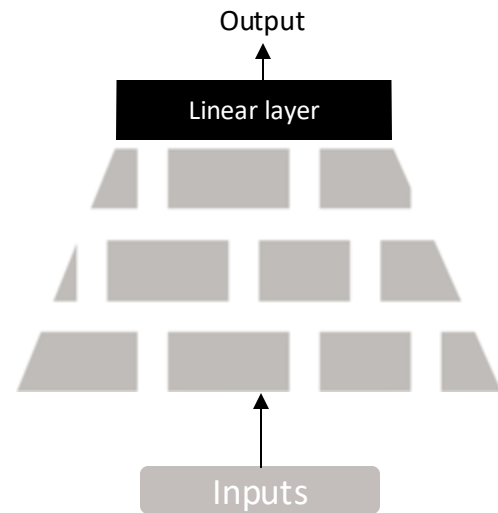
# Hypothesis

Foundation model representations are robust.

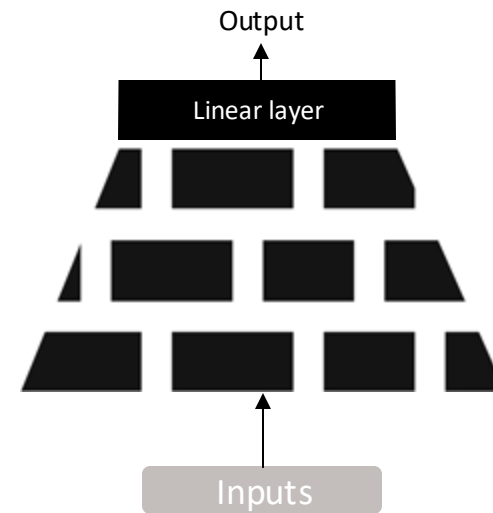We should be able to train on them directly.

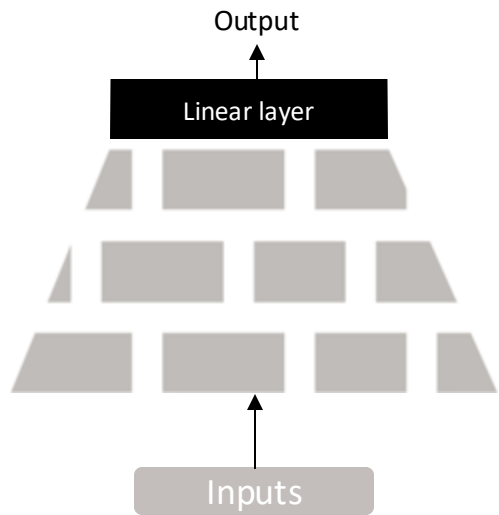# Transfer learning

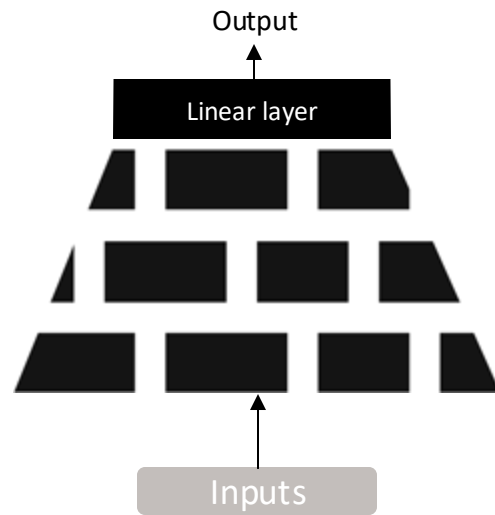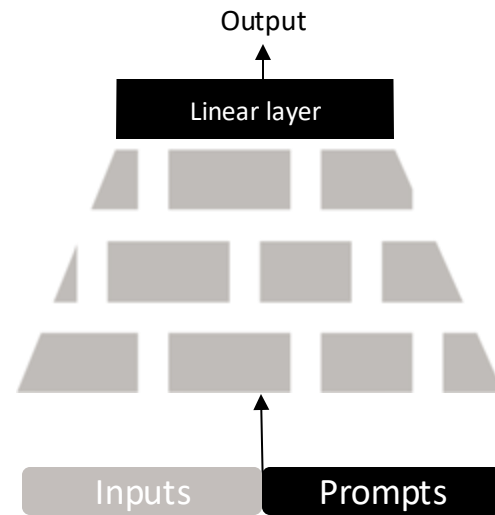No diabetes / Not diabetes

# Transfer learning
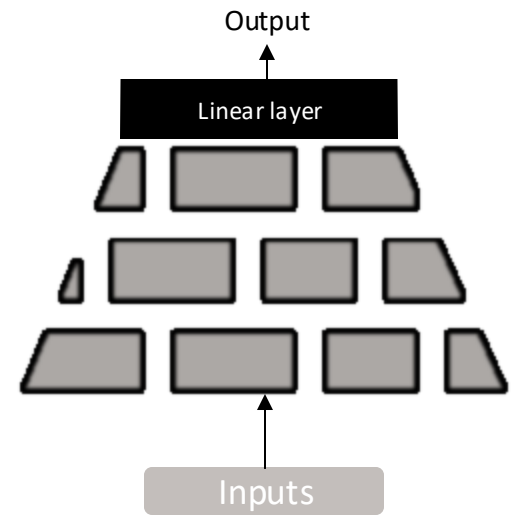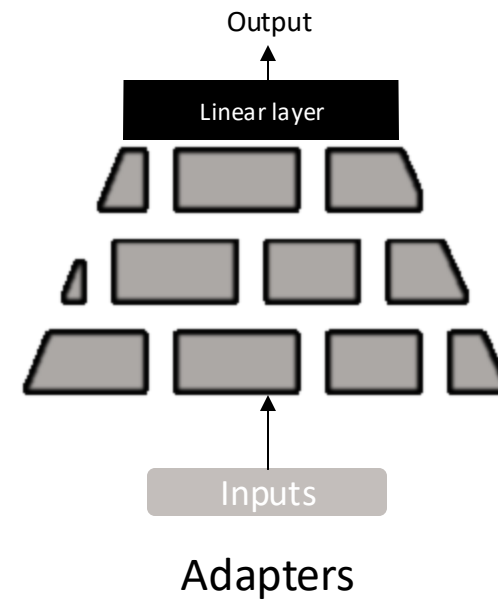


Linear probing

Fine tuning
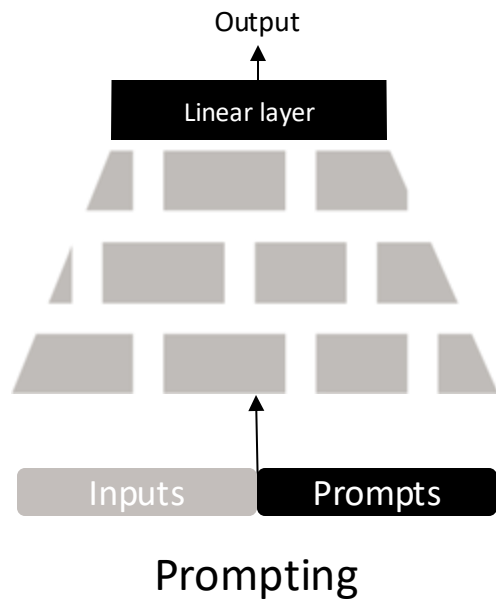
# Efficient adaptation
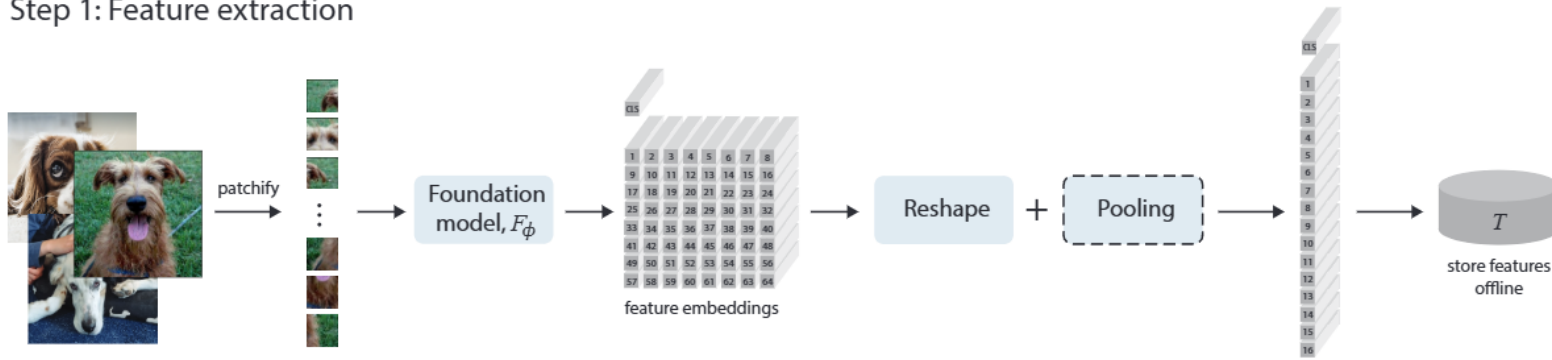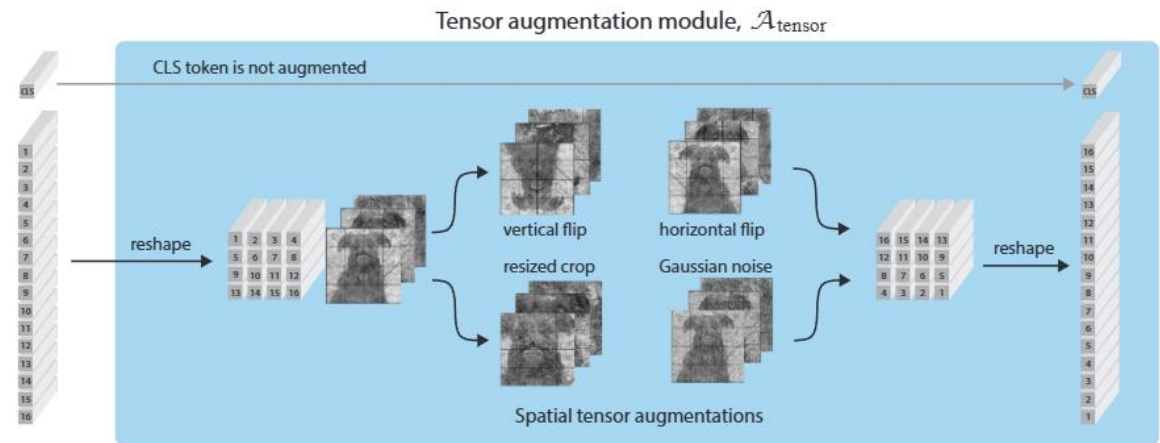


Linear probing

Fine tuning

Prompting

Adapters

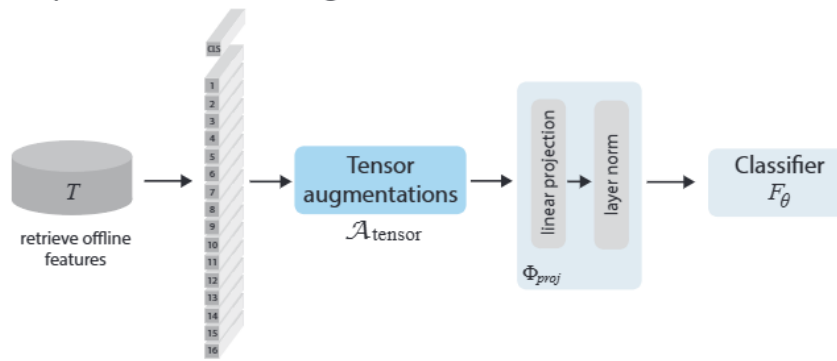# Motivation



Prompting

Adapters

Foundation models are still there.

Can we decouple them?

# LOFFTA



Step 1: Feature extraction

Step 2: Classifier training

# Results

| | Method | Size | APTOS, $\kappa \uparrow$ $n = 3{,}662$ | AID, Acc. $\uparrow$ $n = 10{,}000$ | DDSM, AUC $\uparrow$ $n = 10{,}239$ | ISIC, Rec. $\uparrow$ $n = 25{,}333$ | NABirds, Acc. $\uparrow$ $n = 48{,}562$ | TP, Im/sec $\uparrow$ Train (Infer.) | Mem.,Gb $\downarrow$ Training |
|---|---|---|---|---|---|---|---|---|---|
| ViT-B | Frozen + linear | 256 | $88.6 \pm 0.3.$ | $90.9 \pm 0.1$ | $90.3 \pm 0.2$ | $51.7 \pm 1.0$ | $86.0 \pm 0.1$ | 153 (313) | **1.8** |
| | LOFF | 256 | $89.6 \pm 0.2$ | $91.9 \pm 0.3$ | $94.2 \pm 1.2$ | $70.8 \pm 2.1$ | $83.0 \pm 0.1$ | **228** (236) | 13.2 |
| | LOFF-TA | | $90.4 \pm 0.6$ | $92.3 \pm 0.7$ | $94.4 \pm 0.1$ | $72.8 \pm 1.7$ | $83.5 \pm 0.3$ | 227 (236) | 13.2 |
| | LOFF + Pool | 512 | $89.4 \pm 1.5.$ | $93.2 \pm 0.6$ | $95.3 \pm 0.5$ | $74.3 \pm 1.5$ | $86.2 \pm 0.3$ | **228** (61) | 13.2 |
| | LOFF-TA + Pool | | $\mathbf{90.5 \pm 1.0}$ | $\mathbf{93.7 \pm 0.3}$ | $\mathbf{95.5 \pm 0.1}$ | $\mathbf{77.4 \pm 0.0}$ | $\mathbf{86.8 \pm 0.4}$ | 227 (61) | 13.2 |
| | Unfrozen + linear | 256 | $90.5 \pm 0.9$ | $93.7 \pm 0.8$ | $93.3 \pm 0.9$ | $76.8 \pm 0.7$ | $85.8 \pm 0.1$ | 77 (313) | 28.2 |
| ViT-G | Frozen + linear | 256 | $88.2 \pm 0.3$ | $92.8 \pm 0.2$ | $90.8 \pm 0.6$ | $66.4 \pm 1.1$ | $89.8 \pm 0.2$ | 14 (28) | **7.2** |
| | LOFF | 256 | $88.6 \pm 1.5$ | $93.3 \pm 0.5$ | $94.8 \pm 1.6$ | $73.1 \pm 0.5$ | $87.4 \pm 0.2$ | **222** (27) | 13.2 |
| | LOFF-TA | | $89.9 \pm 0.4$ | $94.0 \pm 0.2$ | $95.3 \pm 0.1$ | $76.0 \pm 0.7$ | $88.5 \pm 0.2$ | 218 (27) | 13.2 |
| | LOFF + Pool | 512 | $90.3 \pm 0.6$ | $94.1 \pm 0.2$ | $95.4 \pm 0.4$ | $74.0 \pm 1.6$ | $88.8 \pm 0.1$ | **222** (7) | 13.2 |
| | LOFF-TA + Pool | | $\mathbf{91.8 \pm 0.3}$ | $\mathbf{94.6 \pm 0.2}$ | $\mathbf{96.3 \pm 0.6}$ | $\mathbf{79.9 \pm 0.2}$ | $\mathbf{90.1 \pm 0.2}$ | 218 (7) | 13.2 |
| | Unfrozen + linear | 256 | $89.6 \pm 0.6$ | $96.2 \pm 0.1$ | $96.7 \pm 0.2$ | $87.3 \pm 1.3$ | $90.2 \pm 0.1$ | 6 (28) | 345.2 |

# Results

| Method | APTOS, $\kappa$ ↑ $n = 3{,}662$ | AID, Acc. ↑ $n = 10{,}000$ | DDSM, AUC ↑ $n = 10{,}239$ | ISIC, Rec. ↑ $n = 25{,}333$ | NABirds, Acc. ↑ $n = 48{,}562$ |
|---|---|---|---|---|---|
| LOFF-TA | 90.4 ± 0.6 | 92.3 ± 0.7 | 94.4 ± 0.1 | 72.8 ± 1.7 | 83.5 ± 0.3 |
| VPT [20] | 89.6 ± 0.1 | 93.0 ± 0.1 | 91.4 ± 0.3 | 75.2 ± 1.1 | 85.8 ± 0.2 |
| VPT + LOFF-TA | 90.8 ± 0.4 | 93.1 ± 0.3 | 92.4 ± 0.3 | 79.7 ± 0.9 | 83.7 ± 0.1 |
| SSF [31] | 90.2 ± 0.1 | 92.1 ± 0.2 | 96.7 ± 0.6 | 76.4 ± 0.9 | 88.2 ± 0.0 |
| SSF + LOFF-TA | 91.1 ± 0.7 | 93.1 ± 0.0 | 97.2 ± 0.3 | 81.6 ± 1.5 | 85.6 ± 0.1 |
| AdaptFormer [6] | 89.6 ± 0.6 | 94.3 ± 0.1 | 91.8 ± 0.8 | 82.6 ± 1.0 | 87.1 ± 0.3 |
| AdaptFormer + LOFF-TA | 90.0 ± 0.3 | 94.3 ± 0.2 | 93.2 ± 0.5 | 83.5 ± 0.3 | 85.3 ± 0.2 |

# Limitations

- Slower during inference

- LOFFTA is competitive but not consistently better in performance

# Conclusions

- Foundation models as fixed feature extractors
- Spatial tensor augmentations

# Thank you!

ekonuk@kth.se