# Fantastic Features and Where to Find Them:
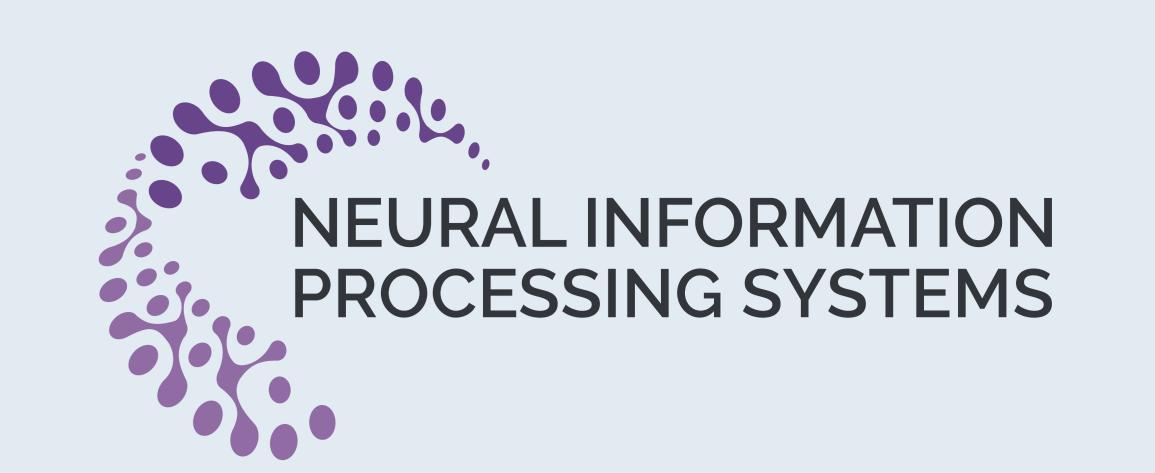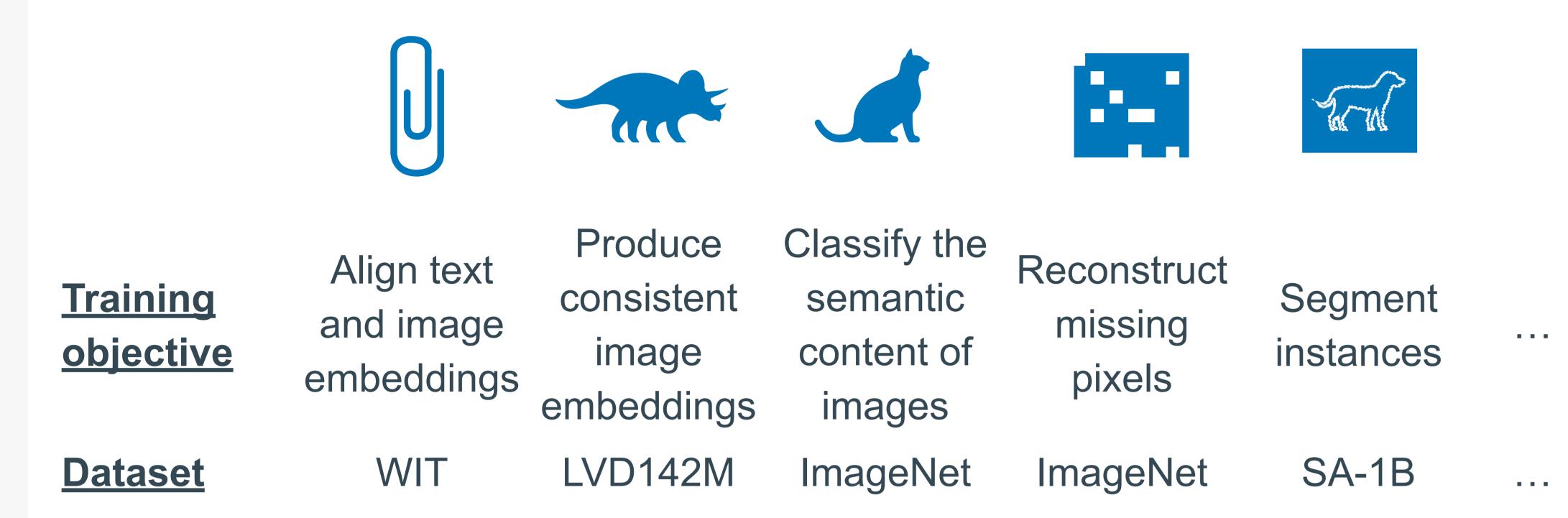## A Probing Method to combine Features from Multiple Foundation Models

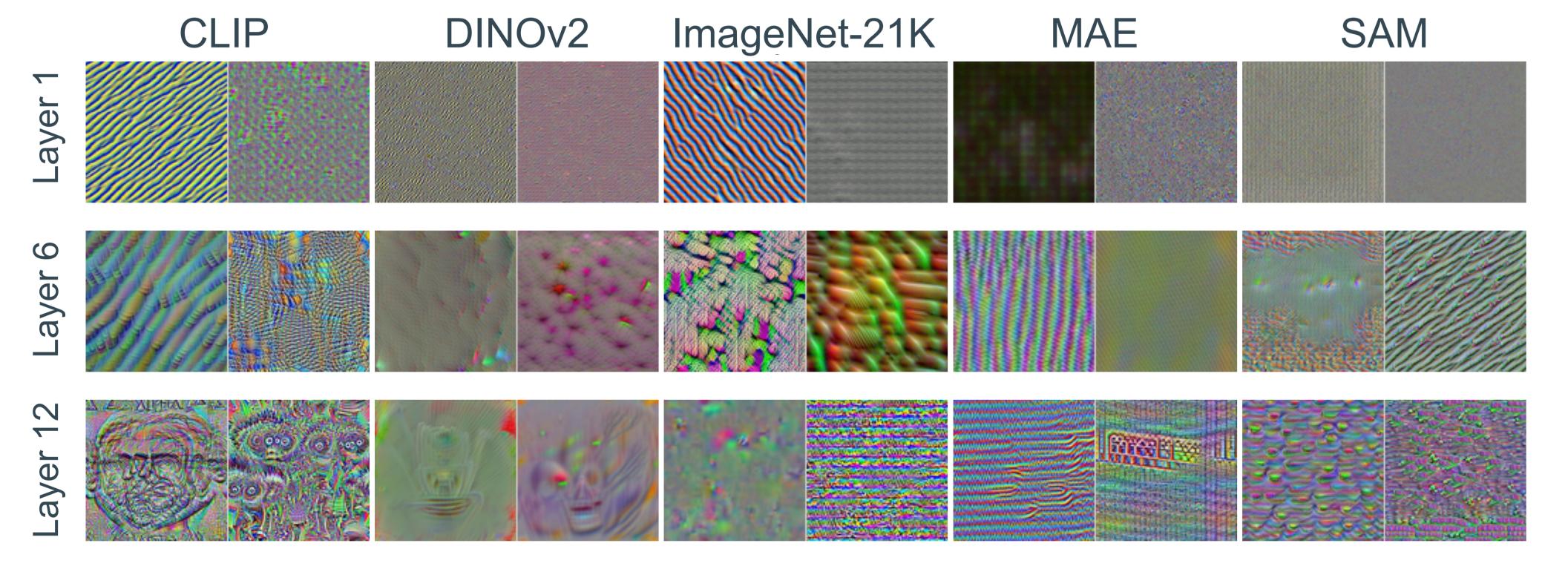*Benjamin Ramtoula, Pierre-Yves Lajoie, Paul Newman, Daniele De Martini*

MRG MOBILE ROBOTICS GROUP OXFORD ROBOTICS INSTITUTE — UNIVERSITY OF OXFORD — POLYTECHNIQUE MONTRÉAL — NEURAL INFORMATION PROCESSING SYSTEMS

## Different foundation models learn different representations

### We now have access to different pre-trained models

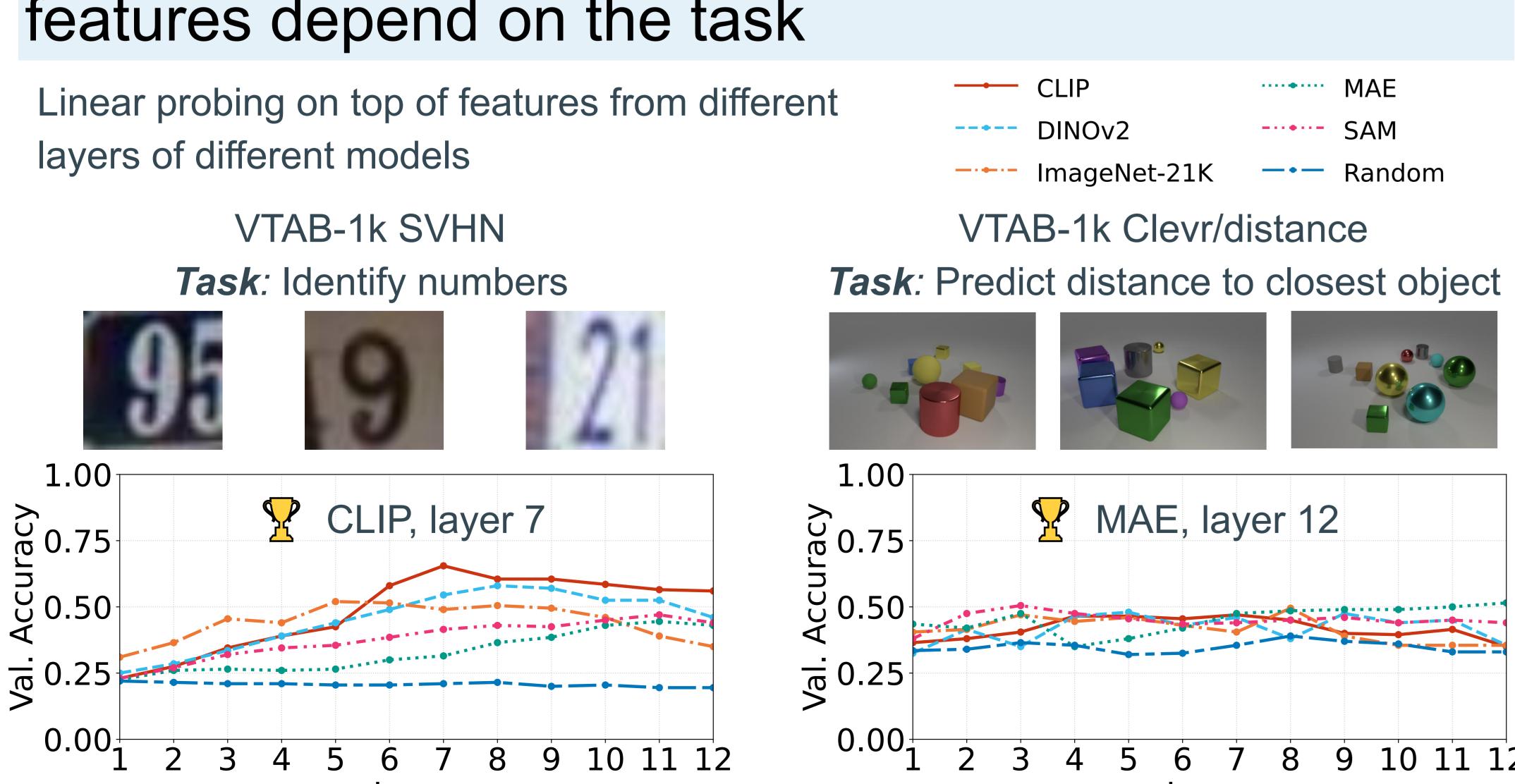| Model | CLIP | DINOv2 | ImageNet | MAE | SAM | ... |
|---|---|---|---|---|---|---|
| Training objective | Align text and image embeddings | Produce consistent image embeddings | Classify the semantic content of images | Reconstruct missing pixels | Segment instances | ... |
| Dataset | WIT | LVD142M | ImageNet | ImageNet | SA-1B | ... |

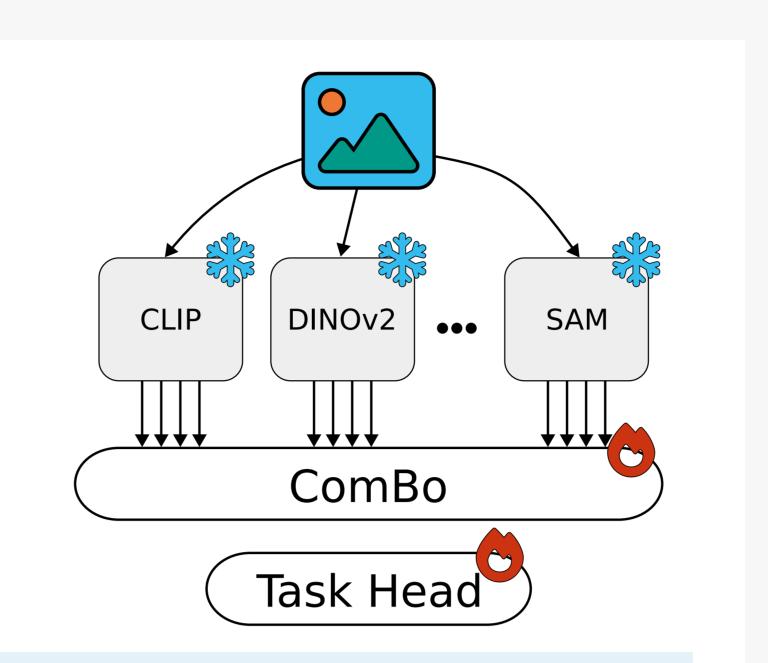### Supervision and data differences affect the representations learned throughout their layers



Images that maximise activations of different neurons

### The model and layer producing the most relevant features depend on the task
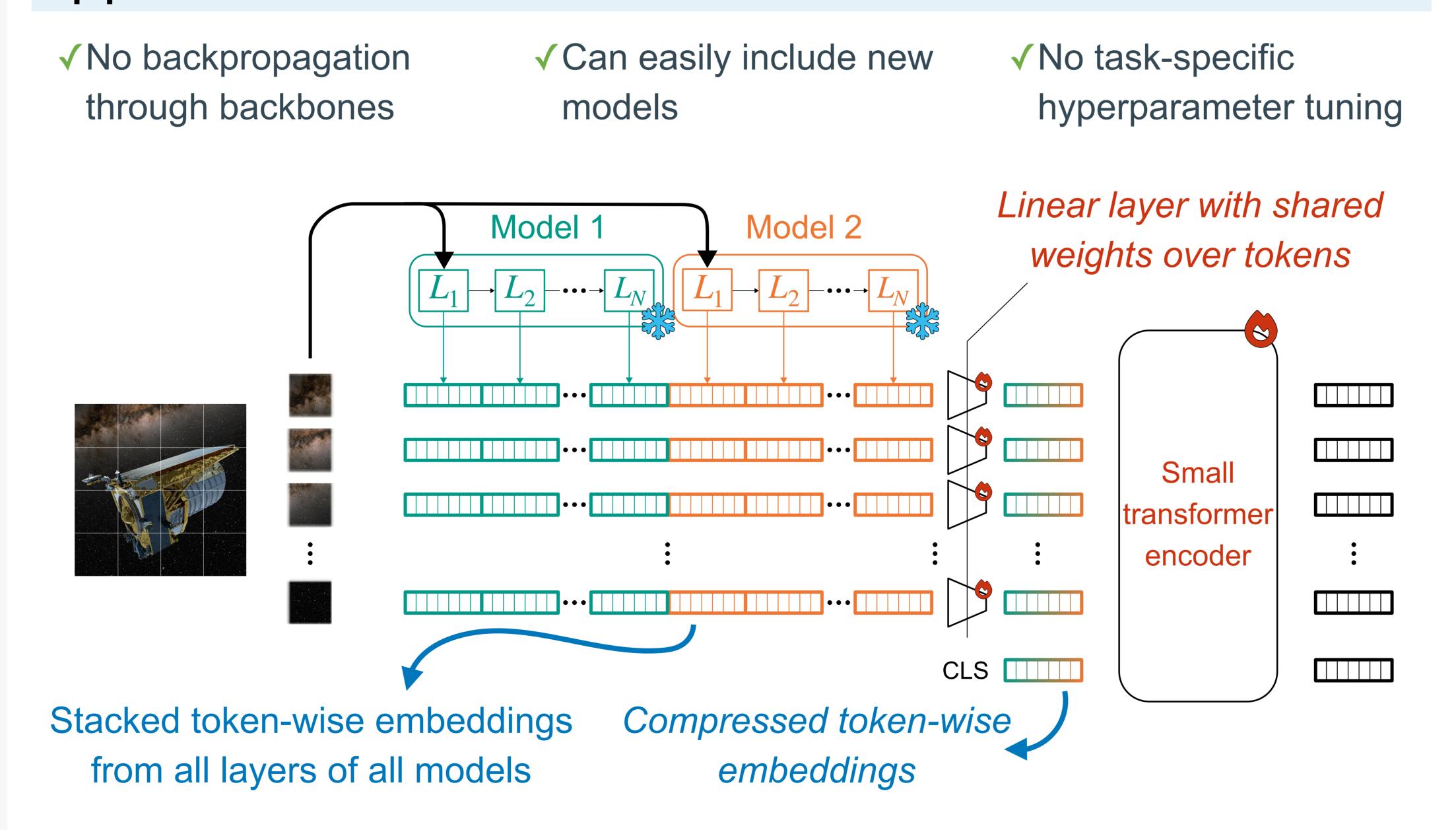
Linear probing on top of features from different layers of different models

Legend: CLIP, DINOv2, ImageNet-21K, MAE, SAM, Random



VTAB-1k SVHN
**Task**: Identify numbers
🏆 CLIP, layer 7

VTAB-1k Clevr/distance
**Task**: Predict distance to closest object
🏆 MAE, layer 12

## We propose an architecture for efficient multi-layer, multi-model feature probing

CLIP  DINOv2  ...  SAM → ComBo → Task Head

### Existing solutions have limitations

| Existing ways to adapt pre-trained models | Examples | Scales to multiple backbones? | Can easily use new backbones? | Can easily be adapted to a new task? |
|---|---|---|---|---|
| Fine-tuning-based approaches | LoRA, Adapter+ | ✗ | ✓ | ✓ |
| Multi-layer probing of frozen features | Head2Toe, SMP | ✓ | ✓ | ✗ |
| Distillation + adaptation | RADIOv2.5 + Adapter+ | ✓ | ✗ | ✓ |

### We address them with **ComBo**, our probing approach to **Com**bine back**Bo**nes

✓ No backpropagation through backbones
✓ Can easily include new models
✓ No task-specific hyperparameter tuning



Model 1  Model 2  *Linear layer with shared weights over tokens*  Small transformer encoder

CLS

*Stacked token-wise embeddings from all layers of all models*
*Compressed token-wise embeddings*

### We can also use ComBo to identify and keep only the most task-relevant models

Using the norm of learned linear layer weights associated with each model to measure their importance:
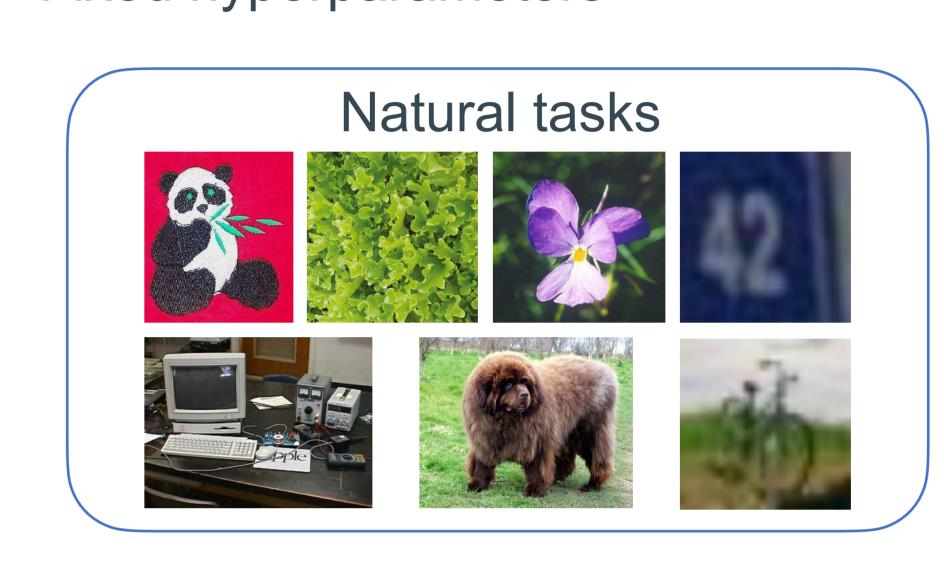1. Train ComBo while minimising each model's importance
2. Inspect weights to measure task-relevance
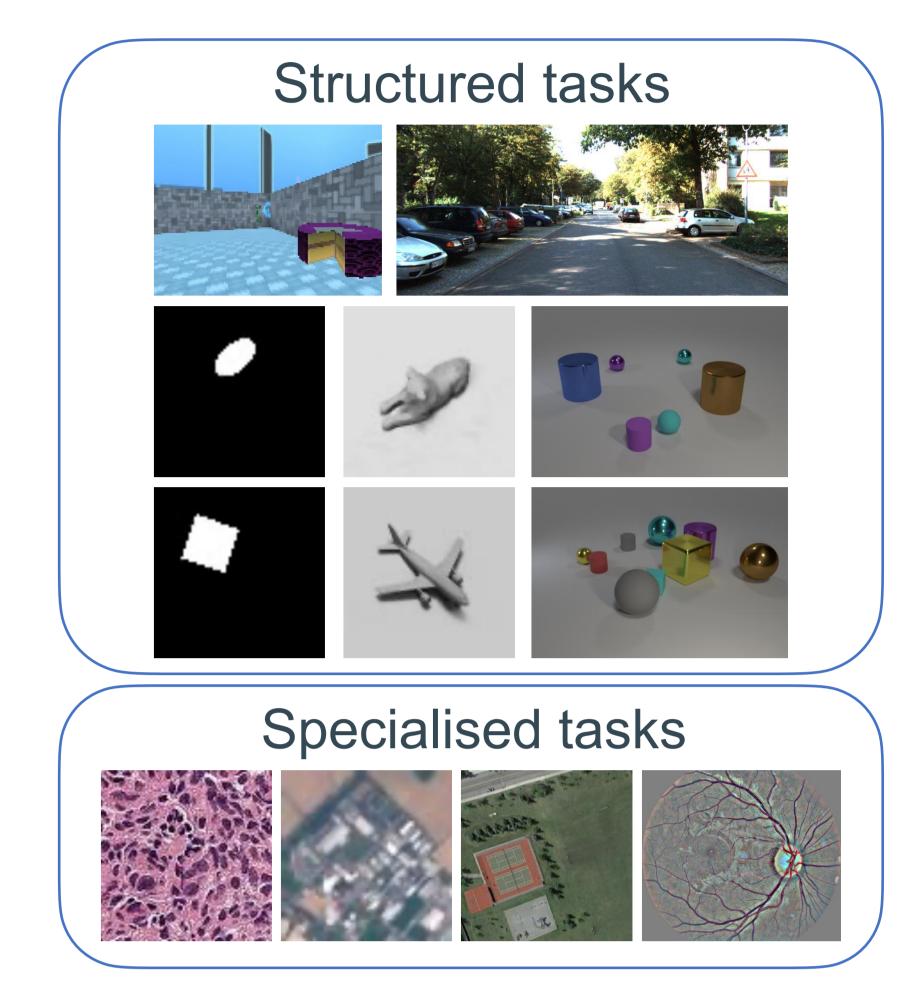3. Retrain using only the most relevant backbones

DFN CLIP
DINOv2
SAM
SigLIP

importance

VTAB-1k tasks

## Why is this useful?
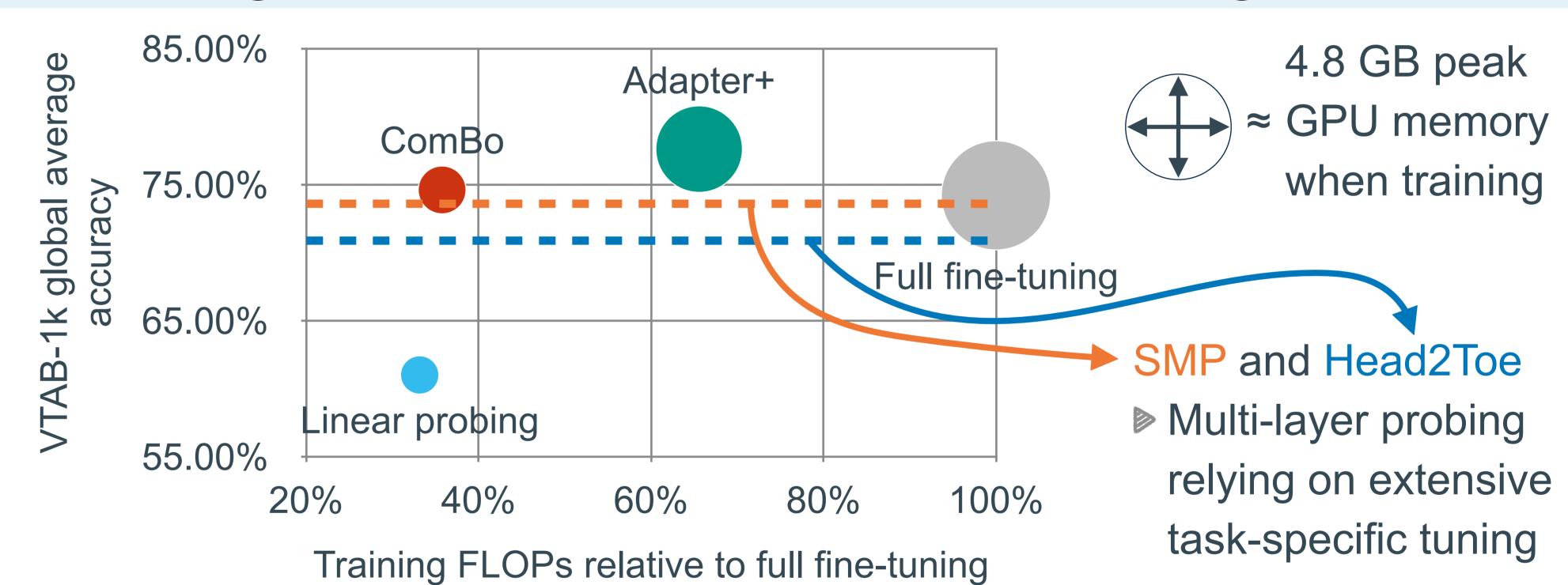
### Experimental setting

- VTAB-1k benchmark:
  - 19 tasks framed as classification
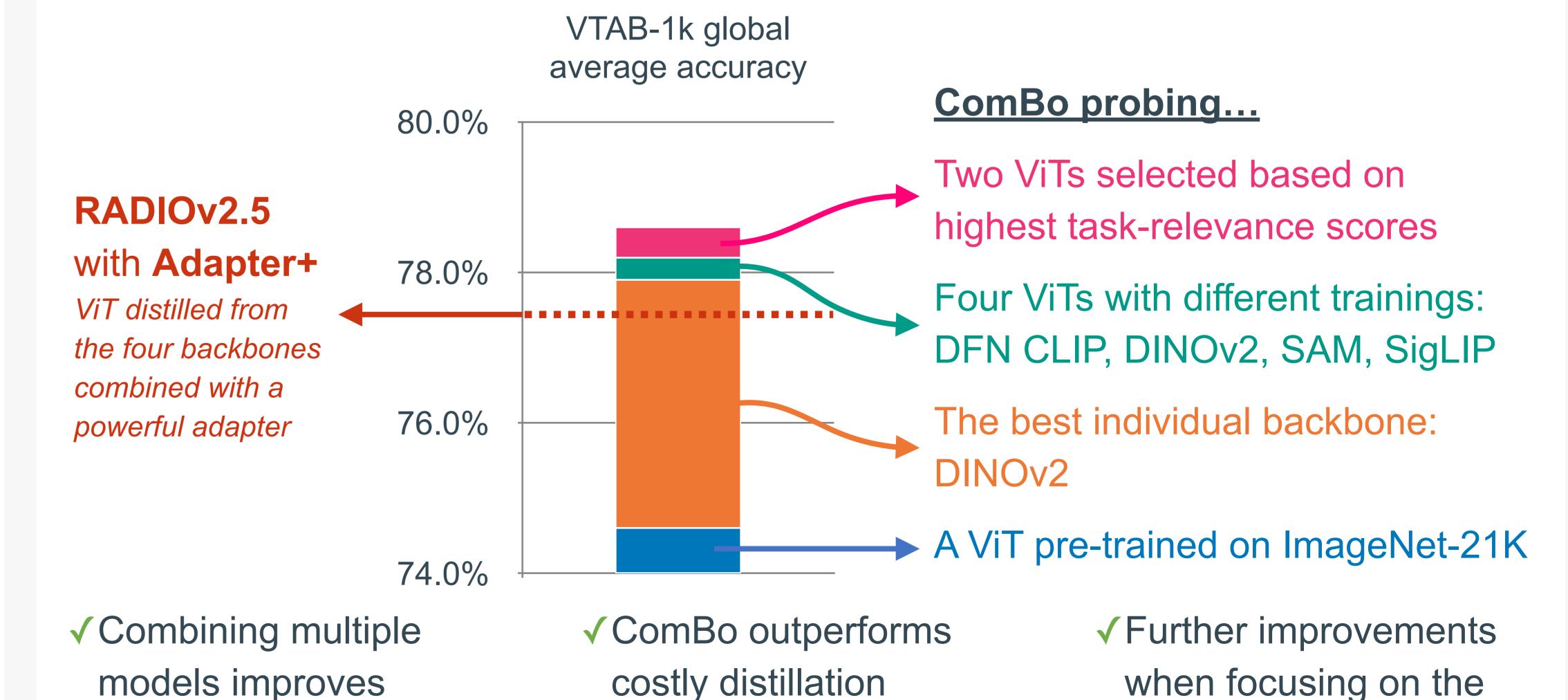  - Only 1000 training images per task
  - Fixed hyperparameters

Natural tasks

Structured tasks

Specialised tasks

### Adapting a ViT-B/16 pre-trained on ImageNet-21K



ComBo — Adapter+ — Linear probing — Full fine-tuning — SMP and Head2Toe

4.8 GB peak ≈ GPU memory when training

▸ Multi-layer probing relying on extensive task-specific tuning

✓ Good performance and minimal compute (5min to train on an RTX 3090 Ti GPU)
  ▸ Enables scaling to multiple backbones
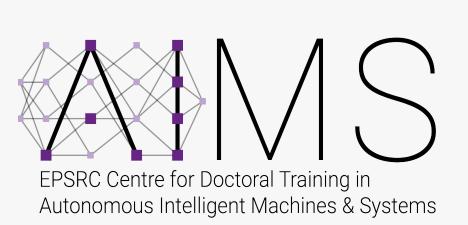
### Probing multiple foundation models at once



VTAB-1k global average accuracy

**RADIOv2.5** with **Adapter+**
*ViT distilled from the four backbones combined with a powerful adapter*

**ComBo probing...**
- Two ViTs selected based on highest task-relevance scores
- Four ViTs with different trainings: DFN CLIP, DINOv2, SAM, SigLIP
- The best individual backbone: DINOv2
- A ViT pre-trained on ImageNet-21K

✓ Combining multiple models improves performance
✓ ComBo outperforms costly distillation
✓ Further improvements when focusing on the most relevant models