

# **Florence: A New Foundation Model for Computer Vision**

Lu Yuan et al.

Florence Team, Microsoft Cloud and AI



Presenter: Minho Park

# A Wide Range of Downstream Tasks

- Achieve state-of-the-art in various benchmarks.

Results from the Paper

✎ Edit

<div>  Ranked #1 on Action Recognition In Videos on Kinetics-400 </div> <div>  Get a GitHub badge </div>									
Task	Dataset	Model	Metric Name	Metric Value	Global Rank	Uses Extra Training Data	Result	Benchmark	
Cross-Modal Retrieval	COCO 2014	Florence	Image-to-text R@1	81.8	# 1	✓	📄	Compare	
			Image-to-text R@5	95.2	# 1	✓	📄	Compare	
			Text-to-image R@1	63.2	# 1	✓	📄	Compare	
			Text-to-image R@5	85.7	# 1	✓	📄	Compare	
Object Detection	COCO minival	Florence-CoSwin-H	box AP	62	# 1	×	📄	Compare	
Object Detection	COCO test-dev	Florence-CoSwin-H	box AP	62.4	# 2	✓	📄	Compare	
Image Classification	ImageNet	Florence-CoSwin-H	Top 1 Accuracy	90.05%	# 7	✓	📄	Compare	
			Top 5 Accuracy	99.02%	# 1	✓	📄	Compare	
Zero-Shot Transfer Image Classification	ImageNet	Florence-CoSwin-H (@384pix)	Accuracy (Private)	83.7	# 3	×	📄	Compare	
Action Recognition In Videos	Kinetics-400	Florence	Top-1 Accuracy	86.8	# 1	×	📄	Compare	
			Top-5 Accuracy	97.5	# 1	×	📄	Compare	

Action Classification	Kinetics-400	Florence (curated FLD-900M pretrain)	Vid acc@1	86.8	# 4	✓	📄	Compare	
			Vid acc@5	97.5	# 2	✓	📄	Compare	
Action Classification	Kinetics-600	Florence (curated FLD-900M pretrain)	Top-1 Accuracy	88.0	# 3	✓	📄	Compare	
			Top-5 Accuracy	97.9	# 3	✓	📄	Compare	
Action Recognition In Videos	Kinetics-600	Florence	Top-1 Accuracy	88.0	# 1	×	📄	Compare	
			Top-5 Accuracy	97.9	# 1	×	📄	Compare	
Video Retrieval	MSR-VTT-1kA	Florence	text-to-video R@1	37.6	# 8	✓	📄	Compare	
			text-to-video R@5	63.8	# 8	✓	📄	Compare	
			text-to-video R@10	72.6	# 8	✓	📄	Compare	
Visual Question Answering	VQA v2 test-dev	Florence	Accuracy	80.16	# 1	×	📄	Compare	
Visual Question Answering	VQA v2 test-std	Florence	overall	80.36	# 2	×	📄	Compare	

Benchmarks in 2022.01.14

# Contribution

---

- A new paradigm of building a computer vision foundation model, Florence.
- Florence successfully extends to different tasks along space, time, and modality, with great transferability.
- Achieves new SOTA results on a wide range of vision benchmarks.

# Foundation Models

- Any model that is trained from broad data at scale that is capable of being adapted (e.g., fine-tuned) to a wide range of downstream tasks.
- Foundation models become promising due to their impressive performance and generalization capabilities.

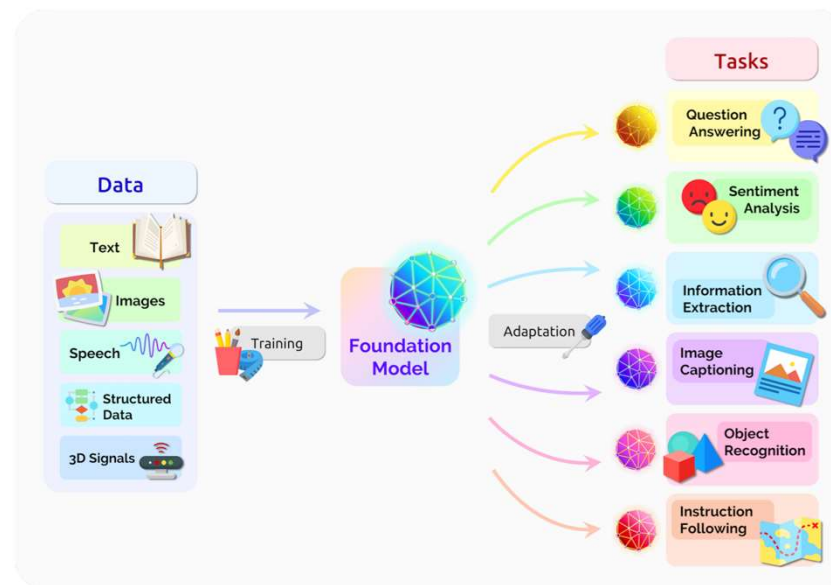
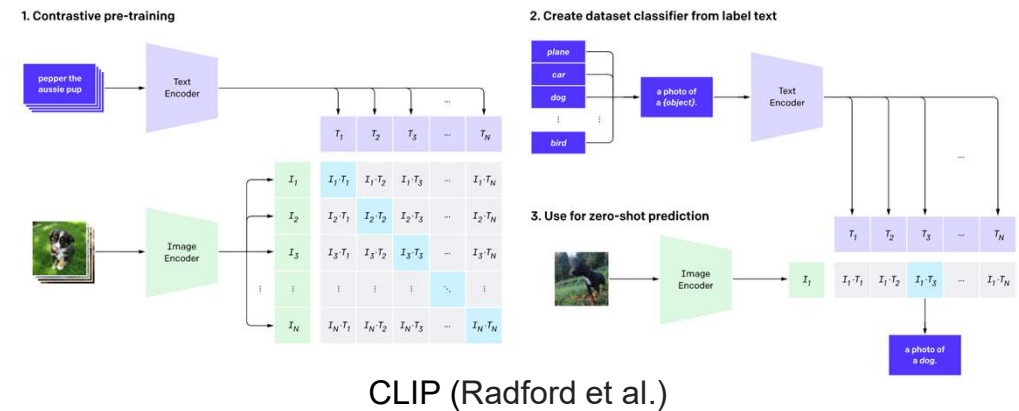
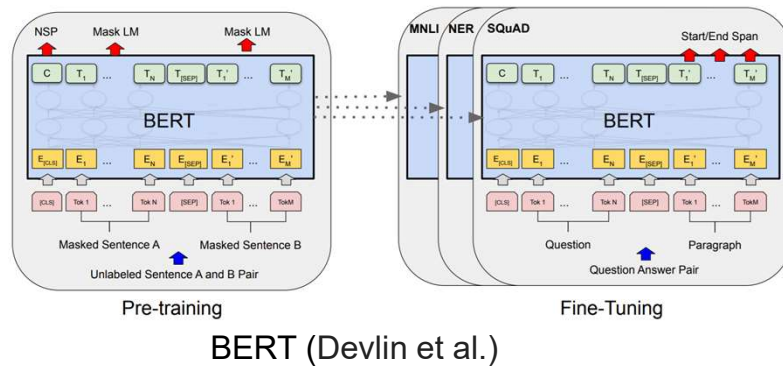


Fig. 2. A foundation model can centralize the information from all the data from various modalities. This one model can then be adapted to a wide range of downstream tasks.

# Foundation Models

- NLP: BERT, GPT
- CV: CLIP, ALIGN, Wu Dao
- However, such models are restricted to image to text mapping.



# What is the foundation model for computer vision?

- Spectrum of tasks in a problem space.
- Space axis: from coarse to fine-grained.
- Time axis: from static to dynamic.
- Modality axis: from RGB only to multiple senses.
- Three axes are orthogonal.

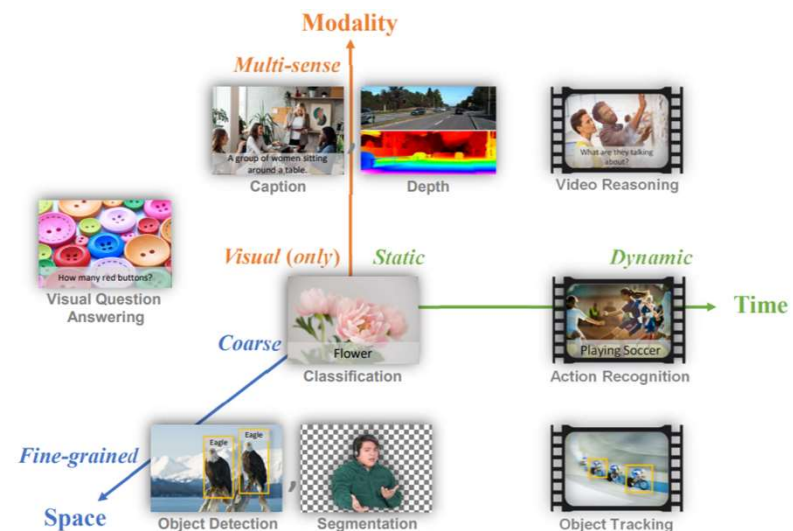


Figure 1. Common computer vision tasks are mapped to a *Space-Time-Modality* space. A computer vision foundation model should serve as general purpose vision system for all of these tasks.

# Dataset Curation

---

- They construct FLD-900M (FLorenceDataset) with 3B image text pairs with selection and post-filtering.  
900M free-form texts (ranging from one word, phrase to sentences)
  - 9.7M unique queries and 7.5B tokens in total.
- Legal and ethical constraints.
- Simple hash-based near-duplicate image removal.
- Small-size image removal.
- Image-text relevance.
- Text should include 500,000 queries (CLIP).  
Occurring at least 100 times in the English version of Wikipedia.
- Class balance the results by including up to 20,000 (image, text) pairs per query (CLIP).  
ImageNet
- Consider informativeness learnability (CLIP).

# Florence

- Florence = Pretrained Models + Adaptation Models

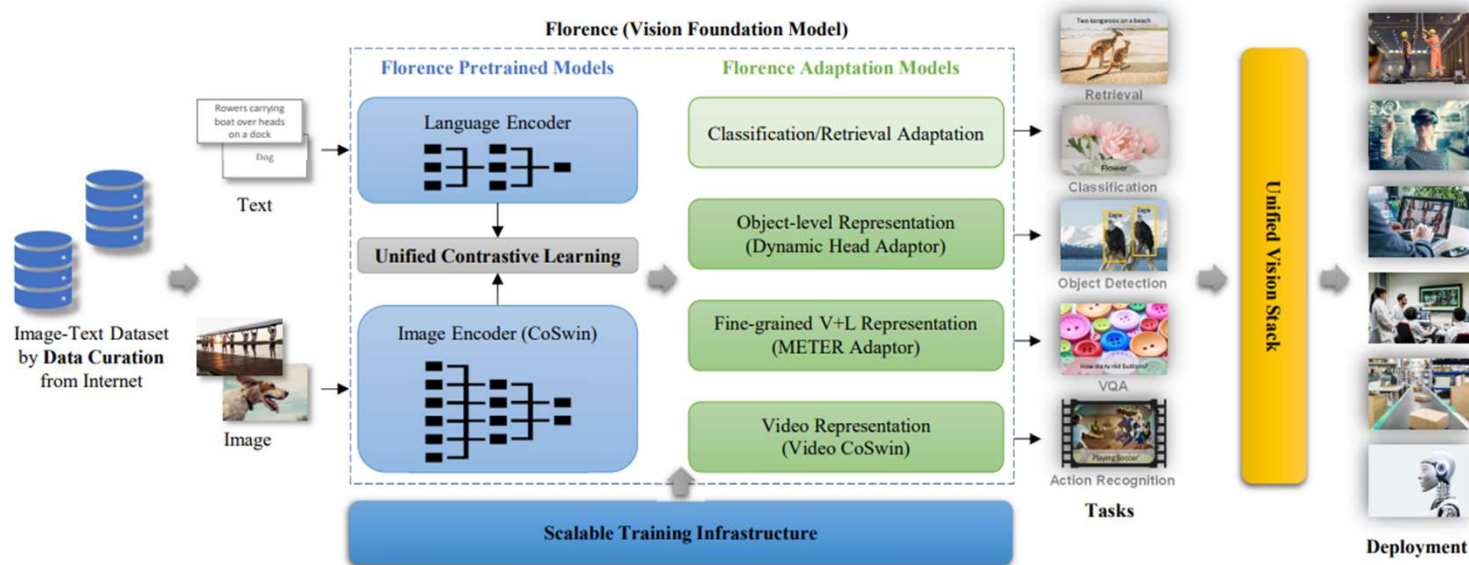


Figure 2. Overview of building Florence. Our workflow consists of data curation, unified learning, Transformer architectures and adaption. It shows the foundation model can be adapted to various downstream tasks and finally integrated into modern computer vision system to power real-world vision and multimedia applications. Compared with existing image-text pretraining models (Radford et al., 2021; Jia et al., 2021; Wud), mainly limited on cross-modal shared representation for classification and retrieval (illustrated by light-green adaptation module), Florence expands the representation to support object level, multiple modality, and videos respectively.



# Unified Image-Text Contrastive Learning (UniCL)

---

- CLIP implicitly assumes that each image-text pair has its unique caption.
- However, in web-scale data, multiple images can be associated with identical captions.
- For example, in FLD-900M, there are 350M image-text pairs where there are more than one images corresponding to one identical text.
- UniCL uses image-text pair as a triplet  $(x, t, y)$ .
  - $x$  is the image,  $t$  is the language description (i.e., hash value), and  $y$  is the language label (i.e., hash key).  
generating propt templates such as “A photo of the [WORD]”
- Thus, all image-text pairs mapped to the same label  $y$  are regarded as positive in our universal image-text contrastive learning. Others are still regarded as negative.

# Unified Image-Text Contrastive Learning (UniCL)

- Given a mini-batch  $\mathcal{B}$ ,  $\mathcal{L} = \mathcal{L}_{i2t} + \mathcal{L}_{t2i}$

$$\mathcal{L}_{i2t} = - \sum_{i \in \mathcal{B}} \frac{1}{|\mathcal{P}(i)|} \sum_{k \in \mathcal{P}(i)} \log \frac{\exp(\tau u_i v_k)}{\sum_{j \in \mathcal{B}} \exp(\tau u_i v_j)}$$

$$\mathcal{L}_{t2i} = - \sum_{j \in \mathcal{B}} \frac{1}{|\mathcal{Q}(j)|} \sum_{k \in \mathcal{Q}(j)} \log \frac{\exp(\tau u_k v_j)}{\sum_{i \in \mathcal{B}} \exp(\tau u_i v_j)}$$

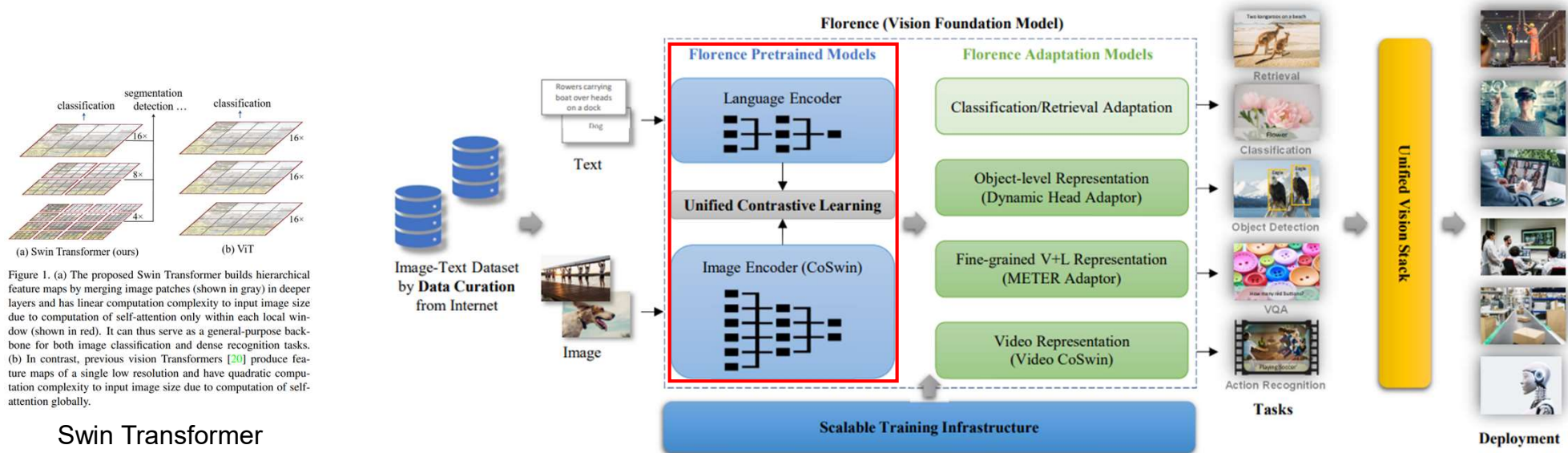
- where  $f_\theta, f_\phi$  are image encoder and text encoder, respectively.
- $u = \frac{f_\theta(x)}{\|f_\theta(x)\|}$ ,  $v = \frac{f_\phi(x)}{\|f_\phi(x)\|}$  are normalized visual feature vector, and  $\tau$  is a learnable temperature.
- $\mathcal{P}(i) = \{k | k \in \mathcal{B}, y_k = y_i\}$ ,  $\mathcal{Q}(j) = \{k | k \in \mathcal{B}, y_k = y_j\}$ .

	$\mathcal{L}_{t2i}$				
	$\mathcal{T}_1$	$\mathcal{T}_2$	$\mathcal{T}_3$	$\dots$	$\mathcal{T}_N$
$\mathcal{I}_1$	1	0	1		0
$\mathcal{I}_2$	0	1	0		0
$\mathcal{I}_3$	1	0	1		0
$\vdots$				$\ddots$	
$\mathcal{I}_N$	0	0	0		1

$\mathcal{L}_{i2t}$

# Transformer-based Florence Pretrained Models

- 12-layer transformer + CoSwin Transformer (CvT)



# Florence Adaptation Models

- There are three additional adaptations (space, time, and modality axes).

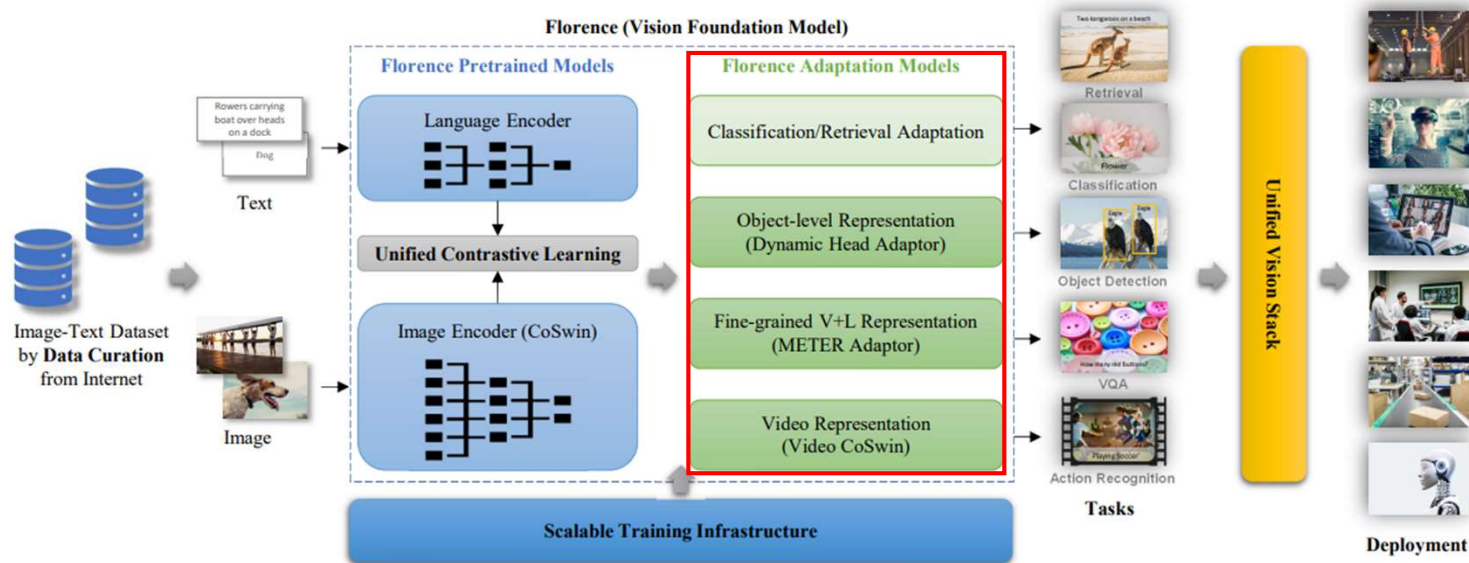


Figure 2. Overview of building Florence. Our workflow consists of data curation, unified learning, Transformer architectures and adaption. It shows the foundation model can be adapted to various downstream tasks and finally integrated into modern computer vision system to power real-world vision and multimedia applications. Compared with existing image-text pretraining models (Radford et al., 2021; Jia et al., 2021; Wud), mainly limited on cross-modal shared representation for classification and retrieval (illustrated by *light-green* adaptation module), Florence expands the representation to support object level, multiple modality, and videos respectively.

Space axis: coarse to fine-grained

# Object-level Visual Representation Learning

- They add an adaptor Dynamic Head (Dai et al., 2021a) (or Dynamic DETR (Dai et al., 2021b)), a unified attention mechanism for the detection head, to the pretrained image encoder (i.e., CoSwin).
- Deploy three attention mechanisms, each on one of the orthogonal dimensions of the tensor, i.e., level-wise, spatial-wise, and channel wise.
- FLOD-9M for object detection pre-training.
  - Merge several well-known object detection datasets, including COCO, LVIS, OpenImages, and Object365.
  - In addition, they generate pseudo bounding boxes on ImageNet-22K

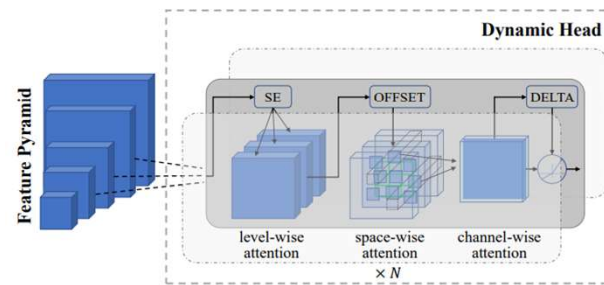


Figure 3. Dynamic Head (Dai et al., 2021a) adapter is used for object-level visual representation learning.

Modality axis: visual only to multi-sense

# Fine-Grained V+L Representation Learning

- They use METER (Dou et al., 2021) adapter to expand to fine-grained vision-language representation (e.g., visual question answering (VQA), image captioning, and fine-grained representation).
- They replace the image encoder of METER with Florence pretrained model CoSwin, and use a pretrained RoBERTa (Liu et al., 2019) as the language encoder.

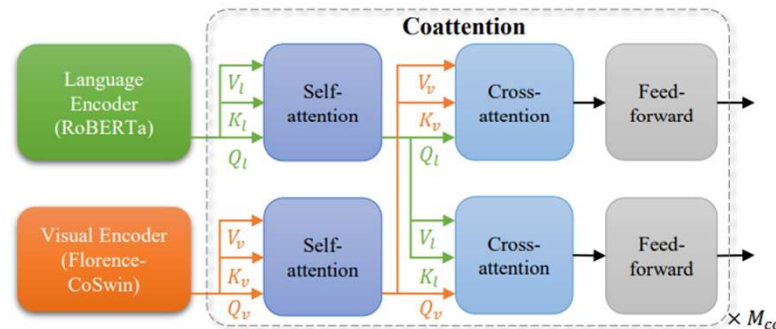


Figure 4. METER (Dou et al., 2021) is used as Florence V+L adaptation model, trained with the image-text matching (ITM) loss and the masked language modeling (MLM) loss.



Time axis: static to dynamic

# Adaptation to Video Recognition

- Video CoSwin: They use 3d convolutional layers instead of 3d patch merging and 2d patch embedding.
- As the initialization to 3D convolutional weights, the pre-trained 2D convolutional weights of CoSwin are duplicated along the temporal dimension and divided by the temporal kernel size to keep the mean and variance of the output unchanged.

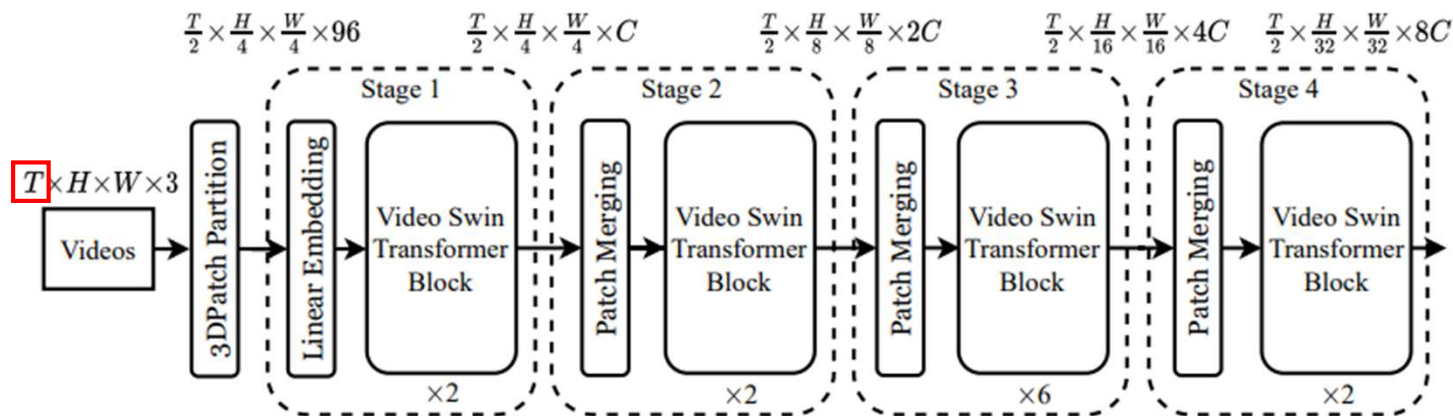


Figure 1: Overall architecture of Video Swin Transformer (tiny version, referred to as Swin-T).

# Experiments

---

- Zero-shot Transfer in Classification
- Linear Probe in Classification
- ImageNet-1K Fine-tune Evaluation
- Few-shot Cross-domain Classification
- Image-Text Retrieval
- Object Detection and Zero-shot Transfer
- V+L Representation Learning
- Zero-shot Text-to-Video Retrieval
- Video Action Recognition



# Zero-shot Transfer in Classification

- Florence can be directly used to predict if an image and a text snippet are semantically matched together in the task dataset (same method of CLIP).

	Food101	CIFAR10	CIFAR100	SUN397	Stanford Cars	FGVC Aircraft	VOC2007	DTD	Oxford Pets	Caltech101	Flowers102	ImageNet
CLIP-ResNet-50x64	91.8	86.8	61.3	48.9	76.0	35.6	83.8	53.4	93.4	90.6	77.3	73.6
CLIP-ViT-L/14 (@336pix)	93.8	<b>95.7</b>	77.5	68.4	78.8	37.2	84.3	55.7	93.5	92.8	78.3	76.2
FLIP-ViT-L/14	92.2	<b>95.7</b>	75.3	73.1	70.8	<b>60.2</b>	-	60.7	92.0	93.0	<b>90.1</b>	78.3
Florence-CoSwin-H (@384pix)	<b>95.1</b>	94.6	<b>77.6</b>	<b>77.0</b>	<b>93.2</b>	55.5	<b>85.5</b>	<b>66.4</b>	<b>95.9</b>	<b>94.7</b>	86.2	<b>83.7</b>

Table 1. Zero-shot transfer of image classification comparisons on 12 datasets: CLIP-ResNet-50x64 (Radford et al., 2021), FLIP-ViT-L/14 (Yao et al., 2021).

# Linear Probe in Classification

- Linear probe as another main metric for evaluating representation quality has been used in most recent studies

	Food101	CIFAR10	CIFAR100	SUN397	Stanford Cars	FGVC Aircraft	VOC2007	DTD	Oxford Pets	Caltech101	Flowers102
SimCLRv2-ResNet-152x3	83.6	96.8	84.5	69.1	68.5	63.1	86.7	80.5	92.6	94.9	96.3
ViT-L/16 (@384pix)	87.4	97.9	89.0	74.9	62.5	52.2	86.1	75.0	92.9	94.7	99.3
EfficientNet-L2 (@800pix)	92.0	<b>98.7</b>	<b>89.0</b>	75.7	75.5	68.4	89.4	82.5	95.6	94.7	97.9
CLIP-ResNet-50x64	94.8	94.1	78.6	81.1	90.5	67.7	88.9	82.0	94.5	95.4	98.9
CLIP-ViT-L/14 (@336pix)	95.9	97.9	87.4	82.2	91.5	71.6	89.9	83.0	95.1	96.0	99.2
Florence-CoSwin-H (@384pix)	<b>96.2</b>	97.6	87.1	<b>84.2</b>	<b>95.7</b>	<b>83.9</b>	<b>90.5</b>	<b>86.0</b>	<b>96.4</b>	<b>96.6</b>	<b>99.7</b>

Table 2. Comparisons of image classification linear probing on 11 datasets with existing state-of-the-art models, including SimCLRv2 (Chen et al., 2020c), ViT (Dosovitskiy et al., 2021a), EfficientNet (Xie et al., 2020), and CLIP (Radford et al., 2021).

# ImageNet-1K Fine-tune Evaluation

---

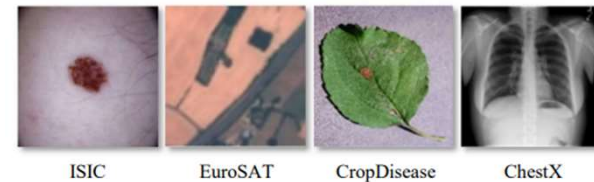
- Florence can be easily adapted to support continual finetuning on target classification tasks.
- State-of-the-art on top-5 accuracy (w/o JFT dataset).

Model	Params	Data	Accuracy	
			Top-1	Top-5
BiT-L-ResNet152x4	928M	300M	87.54	98.46
ALIGN-Efficient-L2	480M	1800M	88.64	98.67
ViT-G/14	1843M	3000M	90.45	-
CoAtNet-7	2440M	3000M	<b><u>90.88</u></b>	-
Florence-CoSwin-H	637M	900M	90.05	<b><u>99.02</u></b>

*Table 3.* Classification fine tuning on ImageNet-1K. Florence is compared with: BiT-L-ResNet152x4 (Kolesnikov et al., 2020), ALIGN-Efficient-L2 (Jia et al., 2021), ViT-G/14 (Zhai et al., 2021), CoAtNet-7 (Dai et al., 2021c) in terms of model scale, data scale and Top-1/Top-5 accuracy.

# Few-shot Cross-domain Classification

- The Cross-Domain Few-Shot learning benchmark (Guo et al., 2020) is used to measure an algorithm's capability to adapt to downstream few-shot target tasks.
- To predict the class, we append a single linear layer as an adapter head to our image encoder CoSwin.
- Training occurs over 100 epochs per episode.



employs ensembles and transductive learning.

	Model	ISIC	EuroSAT	CropD	ChestX	mean
5-shot	CW	57.4	88.1	96.6	29.7	68.0
	Florence	57.1	90.0	97.7	29.3	<b>68.5</b>
20-shot	CW	68.1	94.7	99.2	38.3	75.1
	Florence	72.9	95.8	99.3	37.5	<b>76.4</b>
50-shot	CW	74.1	96.9	99.7	44.4	78.8
	Florence	78.3	97.1	99.6	42.8	<b>79.5</b>

Table 4. Comparison with CW (Liu et al., 2020) (CD-FSL Challenge 2020 Winner) on CD-FSL benchmark. The average result comparison is 74.8 (Florence) vs. 73.9 (CW).

# Image-Text Retrieval

- The zero-shot transfer and fine-tuning performance of Florence for both text and image retrieval.

Method		Flickr30K (1K test set)				MSCOCO (5K test set)			
		Image → Text		Text → Image		Image → Text		Text → Image	
		R@1	R@5	R@1	R@5	R@1	R@5	R@1	R@5
Zero-shot	ImageBERT (Qi et al., 2020)	70.7	90.2	54.3	79.6	44.0	71.2	32.3	59.0
	UNITER (Chen et al., 2020d)	83.6	95.7	68.7	89.2	-	-	-	-
	CLIP (Radford et al., 2021)	88.0	98.7	68.7	90.6	58.4	81.5	37.8	62.4
	ALIGN (Jia et al., 2021)	88.6	98.7	75.7	<u>93.8</u>	58.6	83.0	45.6	69.8
	FLIP (Yao et al., 2021)	89.8	<u>99.2</u>	75.0	93.4	61.3	84.3	45.9	70.6
	Florence	<u>90.9</u>	<u>99.1</u>	<u>76.7</u>	<u>93.6</u>	<u>64.7</u>	<u>85.9</u>	<u>47.2</u>	<u>71.4</u>
Fine-tuned	GPO (Chen et al., 2020a)	88.7	98.9	76.1	94.5	68.1	90.2	52.7	80.2
	UNITER (Chen et al., 2020d)	87.3	98.0	75.6	94.1	65.7	88.6	52.9	79.9
	ERNIE-ViL (Yu et al., 2020)	88.1	98.0	76.7	93.6	-	-	-	-
	VILLA (Gan et al., 2020)	87.9	97.5	76.3	94.2	-	-	-	-
	Oscar (Li et al., 2020)	-	-	-	-	73.5	92.2	57.5	82.8
	ALIGN (Jia et al., 2021)	95.3	99.8	84.9	97.4	77.0	93.5	59.9	83.3
	FLIP (Yao et al., 2021)	96.6	<u>100.0</u>	87.1	97.7	78.9	94.4	61.2	84.3
	Florence	<u>97.2</u>	<u>99.9</u>	<u>87.9</u>	<u>98.1</u>	<u>81.8</u>	<u>95.2</u>	<u>63.2</u>	<u>85.7</u>

Table 5. Image-text retrieval comparisons on Flickr30K and MSCOCO datasets (zero-shot and fine-tuned).

# Object Detection and Zero-shot Transfer

---

- Florence is more desirable for object detection since its adaptation helps learn visual representation at the object level.
- Fine-tuning

Benchmark	Model	AP
<i>COCO miniVal</i>	DyHead	60.3
	Soft Teacher	60.7
	Florence	<b><u>62.0</u></b>
<i>COCO test-Dev</i>	DyHead	60.6
	Soft Teacher	61.3
	Florence	<b><u>62.4</u></b>
<i>Object365</i>	Multi-dataset Detection	33.7
	Florence	<b><u>39.3</u></b>
<i>Visual Genome</i>	VinVL	13.8
	Florence	<b><u>16.2</u></b>

Table 6. Object detection fine tuning comparisons with state-of-the-art methods, including DyHead (Dai et al., 2021a), Soft Teacher (Xu et al., 2021b), Multi-dataset Detection (Zhou et al., 2021), VinVL (Zhang et al., 2021b).



# Object Detection and Zero-shot Transfer

- In our zero-shot transfer setting, object proposal and object classification are decoupled into two tasks.
  - Object proposal discriminates object from background, ignoring semantics of object categories.
  - Classification, on the other hand, focuses on object semantics for each bounding box proposal.
- They freeze the CoSwin backbones and pre-train the Dynamic Head on FLOD-9M by neglecting semantics from each object bounding box.

		Aquarium	BCCD	Chess Pieces	Mask Wearing	Oxford Pets	Packages	Pistols	PKLot	Pothole	Thermal	Wildfire Smoke
	Images	638	364	292	149	3680	26	2986	12416	665	203	737
	Categories	7	3	12	2	37	1	1	2	1	2	1
<i>Fine-tuned</i>	DyHead-Swin-L (full)	53.1	62.6	80.7	52.0	85.9	52.0	74.4	98.0	61.8	75.9	58.7
	DyHead-Swin-L (5-shot)	39.0	40.6	57.3	26.8	47.5	32.8	20.0	22.1	10.8	54.9	14.2
<i>Zero-shot</i>	ZSD	16.0	1.2	0.1	0.6	0.3	58.3	31.5	0.2	2.4	37.4	0.002
	Florence	43.1	15.3	13.4	15.0	68.9	79.6	41.4	31.4	53.3	46.9	48.7

Table 7. Zero-shot transfer in object detection, in comparison with previous state-of-the-art model DyHead (Dai et al., 2021a) (on COCO) fine tuning results on full-set or 5-shot respectively and zero-shot detection baseline model ZSD (Bansal et al., 2018).

# V+L Representation Learning

---

- To evaluate the performance, we fine-tune the pre-trained model on the challenging VQA (Goyal et al., 2017) task

Model	test-dev	test-std
UNITER (Chen et al., 2020d)	73.82	74.02
Visual Parsing (Xue et al., 2021)	74.00	74.17
PixelBERT (Huang et al., 2020)	74.45	74.55
VILLA (Gan et al., 2020)	74.69	74.87
UNIMO (Li et al., 2021c)	75.06	75.27
ALBEF (Li et al., 2021a)	75.84	76.04
VinVL (Zhang et al., 2021b)	76.52	76.60
CLIP-ViL (Shen et al., 2021)	76.48	76.70
METER (Dou et al., 2021)	77.68	77.64
SimVLM (Wang et al., 2021)	80.03	80.34
Florence	<b>80.16</b>	<b>80.36</b>

Table 8. Compare our model with the existing state-of-the-art methods on VQA.



# Zero-shot Text-to-Video Retrieval

- Although Florence is pre-trained on image-text pairs, it can be easily adapted to video tasks.

Method	Pre-training Type	Pre-training Data	R@1	R@5	R@10
MIL-NCE (Miech et al., 2020)	Video	HowTo100M	-	-	32.4
MMV (Alayrac et al., 2020)	Video	HowTo100M, AudioSet	-	-	31.1
VideoCLIP (Xu et al., 2021a)	Video*	HowTo100M	10.4	22.2	30.0
VATT (Akbari et al., 2021)	Video	HowTo100M, AudioSet	-	-	29.7
MCN (Chen et al., 2021)	Image and Video	HowTo100M	-	-	33.8
Frozen-in-Time (Bain et al., 2021)	Image and Video	ImageNet, CC, WebVid-2M	18.7	39.5	51.6
CLIP-ViT-B/16 (Radford et al., 2021)	Image	WIT400M	26.0	49.4	60.7
Florence	Image	FLD-900M	<b><u>37.6</u></b>	<b><u>63.8</u></b>	<b><u>72.6</u></b>

Table 9. Zero-shot text-to-video retrieval results on MSR-VTT 1K-A test set. (\*: Feature extracted from the pre-trained model (Miech et al., 2020), followed by another stage of video-and-language pre-training) The pretraining data used in these existing methods include HowTo100M (Miech et al., 2019), AudioSet (Gemmeke et al., 2017), ImageNet (Deng et al., 2009), CC (Sharma et al., 2018), WebVid-2M (Bain et al., 2021), WIT400M (Radford et al., 2021)

# Video Action Recognition

- They evaluate Florence on fine-tuned video action recognition tasks.

Method	Pretraining Data	Kinetics-400		Kinetics-600		Views	Params
		Top-1	Top-5	Top-1	Top-5		
ViViT-H/16x2	JFT-300M	84.8	95.8	85.8	96.5	$4 \times 3$	648M
VideoSwin-L	ImageNet-22K	84.6	96.5	85.9	97.1	$4 \times 3$	200M
VideoSwin-L	ImageNet-22K	84.9	96.7	86.1	97.3	$10 \times 5$	200M
TokenLearner 16at18+L/10	JFT-300M	85.4	96.3	86.3	97.0	$4 \times 3$	460M
Florence	FLD-900M	<b><u>86.5</u></b>	<b><u>97.3</u></b>	<b><u>87.8</u></b>	<b><u>97.8</u></b>	$4 \times 3$	647M

Table 10. Comparison to state-of-the-art methods, including ViViT (Arnab et al., 2021), VideoSwin (Liu et al., 2021b), TokenLearner (Ryoo et al., 2021), on Kinetics-400 and Kinetics-600. Views indicate  $\#temporal\ clip \times \#spatial\ crop$ .

# Conclusion

---

- A new paradigm of building a computer vision foundation model.
- Although the model size is still below several other existing billion-scale models,
- Florence successfully extends to different tasks along space, time, and modality, with great transferability.
- Achieves new SOTA results on a wide range of vision benchmarks.