

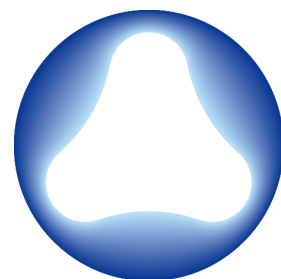
# VideoMAE: Masked Autoencoders are Data-Efficient Learners for Self-Supervised Video Pre-Training

*Zhan Tong<sup>1,2</sup> Yibing Song<sup>2</sup> Jue Wang<sup>2</sup> Limin Wang<sup>1,3</sup>*

*<sup>1</sup> State Key Laboratory for Novel Software Technology, Nanjing University*

*<sup>2</sup> Tencent AI Lab*

*<sup>3</sup> Shanghai AI Lab*



Tencent AI Lab  
腾讯人工智能实验室

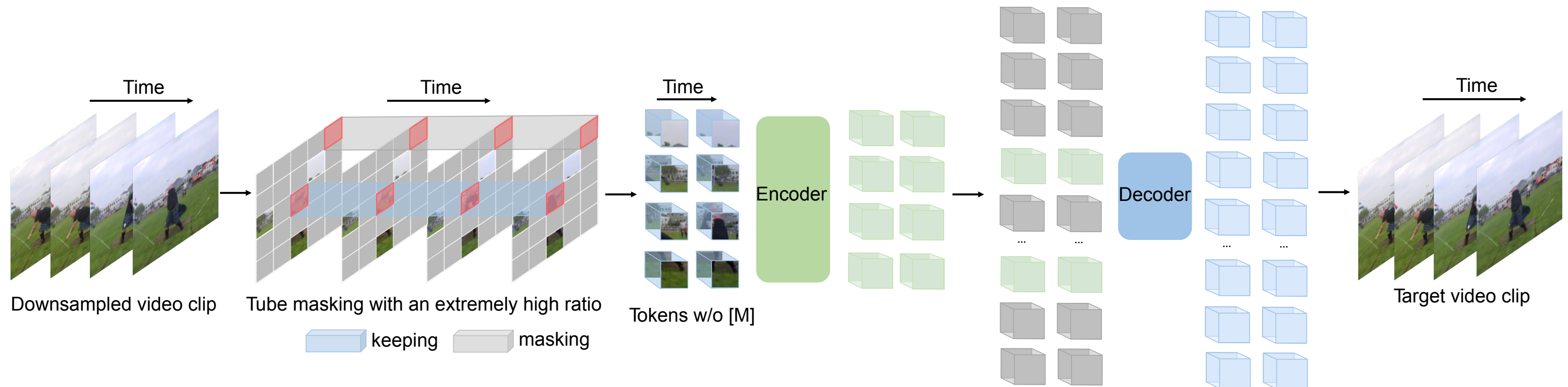


How to **efficiently** train a **vanilla ViT** on the **video** dataset itself  
**without** any pre-trained model or extra data?

# VideoMAE

→ Our VideoMAE attempts to solve it in two aspects

- **Self-supervised** pre-training with **masked autoencoder**
- A new masking strategy: **tube masking** with an **extremely high ratio**



- **and eventually, VideoMAE is**
  - a **simple, data-efficient** method for **self-supervised video pre-training**
    - with **high** performance and **no** extra data **required**

# Experiments

## → Leading performance on Something-Something V2

Method	Backbone	Extra data	Ex. labels	Frames	GFLOPs	Param	Top-1	Top-5
TEINet <sub>En</sub> [39]	ResNet50 <sub>×2</sub>	ImageNet-1K	✓	8+16	99×10×3	50	66.5	N/A
TANet <sub>En</sub> [40]	ResNet50 <sub>×2</sub>		✓	8+16	99×2×3	51	66.0	90.1
TDN <sub>En</sub> [74]	ResNet101 <sub>×2</sub>		✓	8+16	198×1×3	88	69.6	92.2
SlowFast [22]	ResNet101	Kinetics-400	✓	8+32	106×1×3	53	63.1	87.6
MViTv1 [21]	MViTv1-B		✓	64	455×1×3	37	67.7	90.9
TimeSformer [6]	ViT-B	ImageNet-21K	✓	8	196×1×3	121	59.5	N/A
TimeSformer [6]	ViT-L		✓	64	5549×1×3	430	62.4	N/A
ViViT FE [3]	ViT-L	IN-21K+K400	✓	32	995×4×3	N/A	65.9	89.9
Motionformer [50]	ViT-B		✓	16	370×1×3	109	66.5	90.1
Motionformer [50]	ViT-L		✓	32	1185×1×3	382	68.1	91.2
Video Swin [38]	Swin-B		✓	32	321×1×3	88	69.6	92.7
VIMPAC [64]	ViT-L	HowTo100M+DALLE	✗	10	N/A×10×3	307	68.1	N/A
BEVT [76]	Swin-B	IN-1K+K400+DALLE	✗	32	321×1×3	88	70.6	N/A
MaskFeat <sub>↑312</sub> [79]	MViT-L	Kinetics-600	✓	40	2828×1×3	218	75.0	95.0
<b>VideoMAE</b>	ViT-B	Kinetics-400	✗	16	180×2×3	87	69.7	92.3
<b>VideoMAE</b>	ViT-L	Kinetics-400	✗	16	597×2×3	305	74.0	94.6
<b>VideoMAE</b>	ViT-S	<i>no external data</i>	✗	16	57×2×3	22	66.8	90.3
<b>VideoMAE</b>	ViT-B		✗	16	180×2×3	87	70.8	92.4
<b>VideoMAE</b>	ViT-L		✗	16	597×2×3	305	74.3	94.6
<b>VideoMAE</b>	ViT-L		✗	32	1436×1×3	305	<b>75.4</b>	<b>95.2</b>

# Experiments

→ **Leading performance on Kinetics-400**

Method	Backbone	Extra data	Ex. labels	Frames	GFLOPs	Param	Top-1	Top-5
NL I3D [77]	ResNet101	ImageNet-1K	✓	128	$359 \times 10 \times 3$	62	77.3	93.3
TANet [40]	ResNet152		✓	16	$242 \times 4 \times 3$	59	79.3	94.1
TDN <sub>En</sub> [74]	ResNet101		✓	8+16	$198 \times 10 \times 3$	88	79.4	94.4
TimeSformer [6]	ViT-L	ImageNet-21K	✓	96	$8353 \times 1 \times 3$	430	80.7	94.7
ViViT FE [3]	ViT-L		✓	128	$3980 \times 1 \times 3$	N/A	81.7	93.8
Motionformer [50]	ViT-L		✓	32	$1185 \times 10 \times 3$	382	80.2	94.8
Video Swin [38]	Swin-L		✓	32	$604 \times 4 \times 3$	197	83.1	95.9
ViViT FE [3]	ViT-L	JFT-300M	✓	128	$3980 \times 1 \times 3$	N/A	83.5	94.3
ViViT [3]	ViT-H	JFT-300M	✓	32	$3981 \times 4 \times 3$	N/A	84.9	95.8
VIMPAC [64]	ViT-L	HowTo100M+DALLE	✗	10	$N/A \times 10 \times 3$	307	77.4	N/A
BEVT [76]	Swin-B	IN-1K+DALLE	✗	32	$282 \times 4 \times 3$	88	80.6	N/A
MaskFeat <sup>↑352</sup> [79]	MViT-L	Kinetics-600	✗	40	$3790 \times 4 \times 3$	218	87.0	97.4
ip-CSN [68]	ResNet152	<i>no external data</i>	✗	32	$109 \times 10 \times 3$	33	77.8	92.8
SlowFast [22]	R101+NL		✗	16+64	$234 \times 10 \times 3$	60	79.8	93.9
MViTv1 [21]	MViTv1-B		✗	32	$170 \times 5 \times 1$	37	80.2	94.4
MaskFeat [79]	MViT-L		✗	16	$377 \times 10 \times 1$	218	84.3	96.3
<b>VideoMAE</b>	ViT-S	<i>no external data</i>	✗	16	$57 \times 5 \times 3$	22	79.0	93.8
<b>VideoMAE</b>	ViT-B		✗	16	$180 \times 5 \times 3$	87	81.5	95.1
<b>VideoMAE</b>	ViT-L		✗	16	$597 \times 5 \times 3$	305	85.2	96.8
<b>VideoMAE</b>	ViT-H		✗	16	$1192 \times 5 \times 3$	633	<b>86.6</b>	<b>97.1</b>
<b>VideoMAE<sup>↑320</sup></b>	ViT-L	<i>no external data</i>	✗	32	$3958 \times 4 \times 3$	305	86.1	97.3
<b>VideoMAE<sup>↑320</sup></b>	ViT-H		✗	32	$7397 \times 4 \times 3$	633	<b>87.4</b>	<b>97.6</b>

# Experiments

→ **Leading performance on AVA v2.2**

Method	Backbone	Pre-train Dataset	Extra Labels	$T \times \tau$	GFLOPs	Param	mAP
supervised [22]	SlowFast-R101	Kinetics-400	✓	$8 \times 8$	138	53	23.8
CVRL [53]	SlowOnly-R50	Kinetics-400	✗	$32 \times 2$	42	32	16.3
$\rho$ BYOL $_{\rho=3}$ [23]	SlowOnly-R50	Kinetics-400	✗	$8 \times 8$	42	32	23.4
$\rho$ MoCo $_{\rho=3}$ [23]	SlowOnly-R50	Kinetics-400	✗	$8 \times 8$	42	32	20.3
MaskFeat $\uparrow_{312}$ [79]	MViT-L	Kinetics-400	✓	$40 \times 3$	2828	218	37.5
MaskFeat $\uparrow_{312}$ [79]	MViT-L	Kinetics-600	✓	$40 \times 3$	2828	218	38.8
<b>VideoMAE</b>	ViT-S	Kinetics-400	✗	$16 \times 4$	57	22	22.5
<b>VideoMAE</b>	ViT-S	Kinetics-400	✓	$16 \times 4$	57	22	28.4
<b>VideoMAE</b>	ViT-B	Kinetics-400	✗	$16 \times 4$	180	87	26.7
<b>VideoMAE</b>	ViT-B	Kinetics-400	✓	$16 \times 4$	180	87	31.8
<b>VideoMAE</b>	ViT-L	Kinetics-400	✗	$16 \times 4$	597	305	34.3
<b>VideoMAE</b>	ViT-L	Kinetics-400	✓	$16 \times 4$	597	305	37.0
<b>VideoMAE</b>	ViT-H	Kinetics-400	✗	$16 \times 4$	1192	633	<b>36.5</b>
<b>VideoMAE</b>	ViT-H	Kinetics-400	✓	$16 \times 4$	1192	633	<b>39.5</b>
<b>VideoMAE</b>	ViT-L	Kinetics-700	✗	$16 \times 4$	597	305	<b>36.1</b>
<b>VideoMAE</b>	ViT-L	Kinetics-700	✓	$16 \times 4$	597	305	<b>39.3</b>



# Experiments

→ **Leading performance on UCF101 and HMDB51**

Method	Backbone	Extra data	Frames	Param	Modality	UCF101	HMDB51
OPN [35]	VGG	UCF101	N/A	N/A	V	59.6	23.8
VCOP [82]	R(2+1)D	UCF101	N/A	N/A	V	72.4	30.9
CoCLR [29]	S3D-G	UCF101	32	9M	V	81.4	52.1
Vi <sup>2</sup> CLR [18]	S3D	UCF101	32	9M	V	82.8	52.9
<b>VideoMAE</b>	ViT-B	<i>no external data</i>	16	87M	V	<b>91.3</b>	<b>62.6</b>
SpeedNet [5]	S3D-G	Kinetics-400	64	9M	V	81.1	48.8
VTHCL [84]	SlowOnly-R50	Kinetics-400	8	32M	V	82.1	49.2
Pace [73]	R(2+1)D	Kinetics-400	16	15M	V	77.1	36.6
MemDPC [28]	R-2D3D	Kinetics-400	40	32M	V	86.1	54.5
CoCLR [29]	S3D-G	Kinetics-400	32	9M	V	87.9	54.6
RSPNet [12]	S3D-G	Kinetics-400	64	9M	V	93.7	64.7
VideoMoCo [45]	R(2+1)D	Kinetics-400	16	15M	V	78.7	49.2
Vi <sup>2</sup> CLR [18]	S3D	Kinetics-400	32	9M	V	89.1	55.7
CVRL [53]	SlowOnly-R50	Kinetics-400	32	32M	V	92.9	67.9
CVRL [53]	SlowOnly-R50	Kinetics-600	32	32M	V	93.6	69.4
CVRL [53]	Slow-R152 (2×)	Kinetics-600	32	328M	V	94.4	70.6
CORP <sub>f</sub> [32]	SlowOnly-R50	Kinetics-400	32	32M	V	93.5	68.0
$\rho$ SimCLR <sub><math>\rho=2</math></sub> [23]	SlowOnly-R50	Kinetics-400	8	32M	V	88.9	N/A
$\rho$ SwAV <sub><math>\rho=2</math></sub> [23]	SlowOnly-R50	Kinetics-400	8	32M	V	87.3	N/A
$\rho$ MoCo <sub><math>\rho=2</math></sub> [23]	SlowOnly-R50	Kinetics-400	8	32M	V	91.0	N/A
$\rho$ BYOL <sub><math>\rho=2</math></sub> [23]	SlowOnly-R50	Kinetics-400	8	32M	V	92.7	N/A
$\rho$ BYOL <sub><math>\rho=4</math></sub> [23]	SlowOnly-R50	Kinetics-400	8	32M	V	94.2	72.1
<b>VideoMAE(Ours)</b>	ViT-B	Kinetics-400	16	87M	V	<b>96.1</b>	<b>73.3</b>



# VideoMAE: Masked Autoencoders are Data-Efficient Learners for Self-Supervised Video Pre-Training

Code is available at  
<https://github.com/MCG-NJU/VideoMAE>

