044

045

046

047

048

049

050

051

052

053

054

068

069

070

071

072

073

074

075

Snakes and Ladders: Two Steps Up for VideoMamba

Abstract

001 Video understanding requires the extraction of rich spatiotemporal representations, achieved by transformer models 002 through self-attention. Unfortunately, self-attention poses a 003 computational burden. In NLP, Mamba has surfaced as an 004 005 efficient alternative for transformers. However, Mamba's successes do not trivially extend to vision tasks, includ-006 007 ing those in video analysis. In this paper, we theoretically 008 analyze the differences between self-attention and Mamba. We identify two limitations in Mamba's token processing: 009 historical decay and element contradiction. We propose 010 VideoMambaPro (VMP) that addresses these limitations by 011 012 adding masked backward computation and elemental resid-013 ual connections to a VideoMamba backbone. VideoMam-014 baPro models surpass VideoMamba by 1.6-3.0% and 1.1-1.9% top-1 on Kinetics-400 and Something-Something V2, 015 respectively. Even without extensive pre-training, our mod-016 017 els present an attractive and efficient alternative to current 018 transformer models. Moreover, our two solutions are or-019 thogonal to recent advances in Vision Mamba models, and are likely to provide further improvements in future models. 020 021

022 1. Introduction

Video understanding is challenging and requires models 023 024 that can extract rich spatio-temporal representations from video inputs. Transformers are powerful neural networks 025 capable of effectively capturing temporal and spatial in-026 027 formation from videos [19, 29, 50]. Current state-of-the-028 art video understanding models are transformers [41, 51]. 029 At the core of transformers is self-attention [49], the selfalignment between tokens in an input obtained by estimat-030 031 ing the relative importance of a token with respect to all 032 other tokens. The long-range token dependency accounts 033 for much of the success of transformer models [5, 49].

Calculating self-attention is computationally costly,
which eventually limits the application of powerful transformer models in practical settings [17]. Recently, alternative models with lower-cost operators have been proposed

for national language processing (NLP), including S4 [13],038RWKV [38], and RetNet [46]. Among these, Mamba [11]039shows the best performance on long-range and causal tasks040such as language understanding [31] and content-based reasoning [37].042

Motivated by the favorable computational cost, researchers have recently extended Mamba from the NLP domain to the computer vision domain. The core adaptation involved splitting the input image into multiple regions and embedding these as continuous tokens [61]. For video understanding, the recently proposed VideoMamba [22] extracts key frames from videos as the continuous input sequence. However, compared to previous transformer-based methods, VideoMamba's performance on video benchmarks is markedly lower. For example, VideoMamba achieves 82.4% top-1 on Kinetics-400, compared to 85.2% for VideoMAE [47], indicating room for improvement.

In this paper, we first analyze differences in the fea-055 ture extraction capabilities of transformers and Mamba. We 056 identify two limitations of Mamba: historical decay and 057 element contradiction. We then extend VideoMamba [22] 058 to mitigate these limitations. The proposed VideoMam-059 baPro (VMP) addresses historical decay through masked 060 backward computation in the bi-directional Mamba pro-061 cess, allowing the network to better handle historical tokens. 062 We introduce residual connections to Mamba's matrix el-063 ements to tackle element contradiction. VideoMambaPro 064 consistently improves the performance of VideoMamba on 065 video understanding tasks, positioning it as a strong, effi-066 cient competitor to transformers. Our contributions are: 067

- We derive a formal representation of Mamba from the perspective of self-attention and identify two limitations of Mamba in the video analysis domain.
- We propose VideoMambaPro, effectively addressing Mamba's limitations for video understanding.
- We report strong video action recognition performance compared to recent Vision Mamba-based models, and surpass the original VideoMamba by clear margins.

We first discuss related work. Then, we provide a theoretical analysis and introduce VideoMambaPro. Experiments are presented in Section 5. We conclude in Section 6. 078

143

144

145

146

163

164

079 2. Related Work

Transformers. One core aspect of transformers is self-080 attention [49] to achieve long-range interactions by measur-081 ing the similarity between tokens. Self-attention was intro-082 083 duced in the computer vision domain for tasks such as image recognition [26, 45] and object detection [7, 58]. Subse-084 quent works (e.g., [9, 21, 47, 53] extended vision transform-085 ers to the video domain, to achieve superior performance. 086 However, the mechanism of self-attention introduces sig-087 nificant computational overhead because of the similarity 088 089 between all pairs of tokens needs to be calculated.

Alternative models. Recent work has introduced al-090 ternative models with reduced computational complexity, 091 while maintaining the advantages of self-attention [30, 40, 092 59]. SOFT [30] uses Gaussian kernel functions to re-093 094 place the dot-product similarity, which enables a full selfattention matrix to be approximated with a low-rank matrix 095 decomposition. Combiner [40] employs a structured fac-096 097 torization to approximate full self-attention, realizing low computation and memory complexity. 098

Peng et al. [38] propose the Receptance Weighted Key 099 Value (RWKV) architecture that combines self-attention 100 training with an efficient recurrent neural network (RNN) 101 102 inference using a linear attention mechanism. Through parallel computation, a lower, constant-level computational 103 and memory complexity is achieved. RetNet [46] includes 104 another variant of self-attention, by dividing the input into 105 106 multiple chunks. Within each chunk, the self-attention mechanism can be computed in parallel, while information 107 108 is transmitted between chunks based on an RNN.

State-space models. The S4 model completely aban-109 110 dons self-attention and, instead, builds upon a state space model [13]. Instead of calculatings a similarity matrix by 111 performing matrix multiplications for pairs of tokens, it en-112 ables the network to directly learn a global high-order poly-113 114 nomial projection operator (HiPPO) matrix to handle relations between tokens. Additionally, for the simultaneous 115 116 input of multiple tokens, S4 proposes a convolutional processing approach, enabling parallel training and thereby ac-117 celerating the training process. 118

Based on S4, Mamba [11] proposes a selection mech-119 anism where, for each input token, a unique HiPPO ma-120 121 trix [12] is generated. This allows the model to selectively process input tokens, enabling it to focus on or ignore spe-122 123 cific inputs. Due to Mamba's strong representation ability in NLP, and linear-time complexity, it has garnered at-124 125 tention as a promising alternative to transformers. In the 126 computer vision domain, researchers have proposed Vision Mamba [61] and VMamba [27] for tasks such as image clas-127 sification and object detection. 128

In the video domain, several Mamba variants have been
proposed [22, 23, 34]. Their performance is somewhat
lower than expected, with limited understanding of the

causes. We argue that a systematic, mathematical analy-
sis of Mamba from the perspective of self-attention could132reveal shortcomings of Mamba's inner workings. Better
understanding of these limitations allow us to develop im-
provements, and to close the accuracy performance gap
with transformers, while enjoying Mamba's efficiency.132133
134135

3. Theoretical Analysis

We revisit Mamba from the perspective of self-attention.139Then, we analyze its limitations for video understanding.140To address these, we propose VideoMambaPro in Section 4.141

3.1. Mamba from the perspective of self-attention 142

Self-attention. Given an input sequence $X := [x_1, \dots, x_N] \in \mathbb{R}^{N \times D_x}$ of N feature vectors of depth D_x , self-attention [49, 60] computes the output sequence **Y** from **X** following two steps:

Step 1: Similarity matrix computation. The input sequence X is linearly projected onto the three different subspaces query $\mathbf{Q} \in \mathbb{R}^{N \times D}$, key $\mathbf{K} \in \mathbb{R}^{N \times D}$, and value148 $\mathbf{V} \in \mathbb{R}^{N \times D_V}$:150

$$\mathbf{Q} = \boldsymbol{X} \mathbf{W}_Q^{\top}, \mathbf{K} = \boldsymbol{X} \mathbf{W}_K^{\top}, \mathbf{V} = \boldsymbol{X} \mathbf{W}_V^{\top}.$$
(1) 151

with $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{D \times D_x}$, and $\mathbf{W}_V \in \mathbb{R}^{D_v \times D_x}$ 152 the corresponding weight matrices. Specifically, $\mathbf{Q} :=$ 153 $[\boldsymbol{q}_1, \cdots, \boldsymbol{q}_N]^\top$, $\mathbf{K} := [\boldsymbol{k}_1, \cdots, \boldsymbol{k}_N]^\top$, and $\mathbf{V} :=$ 154 $[\boldsymbol{v}_1, \cdots, \boldsymbol{v}_N]^\top$ with vectors $\boldsymbol{q}_i, \boldsymbol{k}_i, \boldsymbol{v}_i$ for $i = 1, \cdots, N$ 155 the query, key, and value vectors, respectively, for input *i*. 156 Based on \mathbf{Q} and \mathbf{K} , similarity matrix $\mathbf{S} \in \mathbb{R}^{N \times N}$ contains the correlations between all query and key vectors, with a softmax function applied to each row of \mathbf{S} : 159

$$\mathbf{S} = \operatorname{softmax}(\mathbf{Q}\mathbf{K}^{\top}/\sqrt{D}). \tag{2}$$
 160

Each component s_{ij} $(i, j = 1, \dots, N)$ represents the similarity score between q_i and k_j . 161

Step 2: Output computation. Output sequence $\mathbf{Y} := [\mathbf{y}_1, \cdots, \mathbf{y}_N]^\top$ is then calculated based on S as:

$$\mathbf{Y} = \mathbf{S}\mathbf{V}.\tag{3}$$

It follows that each output vector y_i $(i = 1, \dots, N)$ can 166 be written in vector form as: 167

$$\boldsymbol{y_i} = \sum_{j=1}^{N} s_{ij} \boldsymbol{v_j}.$$
 (4) 168

Any output vector y_i is a linear combination of vectors 169 $v_j (j = 1, \dots, N)$, with similarity score s_{ij} serving as coefficient. The larger the similarity score, the greater the influence of v_j on output y_i [43]. 172

Mamba. State Space Models (SSMs) serve as the foundation of Mamba [11]. They are based on continuous systems that map 1D functions or sequences, $x(t) \in \mathbb{R}^L \rightarrow$ 175

194

195 196

176 $y(t) \in \mathbb{R}^L$ to output sequences y(t) through a hidden state 177 $h(t) \in \mathbb{R}^N$. Formally, SSM implements the mapping as¹:

$$h(t) = \mathbf{A}h(t-1) + \mathbf{B}x(t), \tag{5}$$

$$y(t) = \mathbf{C}h(t) \tag{6}$$

where $\mathbf{A} \in \mathbb{R}^{N \times N}$ is the evolution matrix of the system, 181 and $\mathbf{B} \in \mathbb{R}^{N \times 1}$, $\mathbf{C} \in \mathbb{R}^{N \times 1}$ are the projection matrices. 182 Often, the input is discrete rather than a continuous func-183 184 tion x(t). Therefore, Mamba performs discretization, ef-185 fectively creating a discrete version of the continuous system. A timescale parameter Δ is used to transform the con-186 tinuous parameters AandB into their discrete counterparts 187 $\overline{\mathbf{A}}, \overline{\mathbf{B}}$, and the transformation typically employs the zero-188 189 order hold method [57]. This process is expressed as:

$$\overline{\mathbf{A}} = \exp(\mathbf{\Delta}\mathbf{A}),\tag{7}$$

$$\overline{\mathbf{B}} = (\mathbf{\Delta}\mathbf{A})^{-1}(\exp(\mathbf{\Delta}\mathbf{A}) - \mathbf{I}) \cdot \mathbf{\Delta}\mathbf{B}, \quad (8)$$

$$h_t = \overline{\mathbf{A}}h_{t-1} + \overline{\mathbf{B}}x_t,\tag{9}$$

$$y_t = \mathbf{C}h_t. \tag{10}$$

Considering that parameters $\overline{\mathbf{A}}, \overline{\mathbf{B}}, \mathbf{C}$ in the original 197 198 SSM are independent of the input data x(t) and cannot be tailored to specific input data, Mamba employs a Selective 199 Scan Mechanism as its core operator. More precisely, three 200 functions $S_B(x), S_C(x), S_{\Delta}(x)$ are introduced to associate 201 parameters $\overline{\mathbf{B}}, \mathbf{C}, and \Delta$ in Eqs. 7–10 to the input data x. 202 Based on $S_{\Delta}(x)$, $\overline{\mathbf{A}}$ can also be associated with the input 203 data x. For example, given the input x_1 , functions $S_{\Delta}(x)$ 204 will produce the corresponding $\overline{\mathbf{A}}_1$ based on Eq. 7, and 205 functions $S_B(x)$ and $S_{\Delta}(x)$ will produce the correspond-206 ing $\overline{\mathbf{B}}_1$ based on Eq. 8. $\overline{\mathbf{C}}_1$ is obtained based on function 207 $S_C(x)$. Following Eqs. 9 and 10, we analyze the process 208 209 to obtain output sequence Y when given an input sequence $\boldsymbol{X} := [\boldsymbol{x_1}, \cdots, \boldsymbol{x_N}] \in \mathbb{R}^{N \times D_x}$ of N feature vectors. Each 210 vector's hidden state is denoted as: 211

$$h_1 = \mathbf{B}_1 \boldsymbol{x}_1, \tag{11}$$

213
$$h_2 = \overline{\mathbf{A}}_2 h_1 + \overline{\mathbf{B}}_2 \boldsymbol{x_2}$$

$$= \overline{\mathbf{A}}_2 \overline{\mathbf{B}}_1 \boldsymbol{x}_1 + \overline{\mathbf{B}}_2 \boldsymbol{x}_2, \qquad (12)$$

215
$$h_3 = \overline{\mathbf{A}}_3 h_2 + \overline{\mathbf{B}}_3 \boldsymbol{x}_3$$

216
$$= \overline{\mathbf{A}}_3 \overline{\mathbf{A}}_2 \overline{\mathbf{B}}_1 \boldsymbol{x}_1 + \overline{\mathbf{A}}_3 \overline{\mathbf{B}}_2 \boldsymbol{x}_2 + \overline{\mathbf{B}}_3 \boldsymbol{x}_3, \quad (13)$$

217

218
$$h_{N} = \overline{\mathbf{A}}_{N} h_{N-1} + \overline{\mathbf{B}}_{N} \boldsymbol{x}_{N}$$
219
$$= \overline{\mathbf{A}}_{N} \overline{\mathbf{A}}_{N-1} \cdots \overline{\mathbf{A}}_{2} \overline{\mathbf{B}}_{1} \boldsymbol{x}_{1} + \overline{\mathbf{A}}_{N} \overline{\mathbf{A}}_{N-1} \cdots \overline{\mathbf{A}}_{3} \overline{\mathbf{B}}_{2} \boldsymbol{x}_{2}$$
220
$$+ \overline{\mathbf{A}}_{N} \overline{\mathbf{B}}_{N-1} \boldsymbol{x}_{N-1} + \overline{\mathbf{B}}_{N} \boldsymbol{x}_{N}.$$
(14)

Eqs. 11–14 can be written in matrix form:

$$\mathbf{H} = \begin{bmatrix} h_1, h_2, h_3, \cdots, h_N \end{bmatrix}^{\top}$$
$$= \begin{bmatrix} \mathbf{\overline{B}}_1 & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{\overline{A}}_2 \mathbf{\overline{B}}_1 & \mathbf{\overline{B}}_2 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{\overline{A}}_3 \mathbf{\overline{A}}_2 \mathbf{\overline{B}}_1 & \mathbf{\overline{A}}_3 \mathbf{\overline{B}}_2 & \mathbf{\overline{B}}_3 & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ (\prod_{j=\mathbf{N}}^2 \mathbf{\overline{A}}_j) \mathbf{\overline{B}}_1 & (\prod_{j=\mathbf{N}}^3 \mathbf{\overline{A}}_j) \mathbf{\overline{B}}_2 & (\prod_{j=\mathbf{N}}^4 \mathbf{\overline{A}}_j) \mathbf{\overline{B}}_3 & \cdots & \mathbf{\overline{B}}_N \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \vdots \\ \mathbf{x}_N \end{bmatrix}.$$
(15) 222

For output sequence $\mathbf{Y} := [\boldsymbol{y}_1, \cdots, \boldsymbol{y}_N]^{\top}$, each vector 223 $\boldsymbol{y}_i \ (i = 1, \cdots, N)$ can be expressed as: 224

$$\boldsymbol{y}_{N} = \overline{\mathbf{C}}_{N} h_{N}, \qquad (16) \qquad \mathbf{225}$$

and in matrix form as:

$$\mathbf{Y} = \begin{bmatrix} \overline{\mathbf{C}}_{1} & 0 & 0 & \cdots & 0\\ 0 & \overline{\mathbf{C}}_{2} & 0 & \cdots & 0\\ 0 & 0 & \overline{\mathbf{C}}_{3} & \cdots & 0\\ \vdots & \vdots & \vdots & \ddots & \vdots\\ 0 & 0 & 0 & \cdots & \overline{\mathbf{C}}_{N} \end{bmatrix} \begin{bmatrix} h_{1}\\ h_{2}\\ h_{3}\\ \vdots\\ h_{N} \end{bmatrix} = \overline{\mathbf{C}}h. \quad (17) \qquad 227$$

By substituting Eq. 15 into Eq. 17, we obtain: 228

$$\mathbf{Y} = \overline{\mathbf{C}} \begin{bmatrix} \overline{\mathbf{B}}_{1} & 0 & 0 & \cdots & 0\\ \overline{\mathbf{A}_{2}\mathbf{B}_{1}} & \overline{\mathbf{B}}_{2} & 0 & \cdots & 0\\ \vdots & \overline{\mathbf{A}_{3}} \overline{\mathbf{A}_{2}} \overline{\mathbf{B}}_{1} & \overline{\mathbf{A}_{3}} \overline{\mathbf{B}}_{2} & \overline{\mathbf{B}}_{3} & \cdots & 0\\ \vdots & \vdots & \vdots & \ddots & \vdots\\ (\prod_{j=\mathbf{N}}^{2} \overline{\mathbf{A}_{j}}) \overline{\mathbf{B}}_{1} & (\prod_{j=\mathbf{N}}^{3} \overline{\mathbf{A}_{j}}) \overline{\mathbf{B}}_{2} & (\prod_{j=\mathbf{N}}^{4} \overline{\mathbf{A}_{j}}) \overline{\mathbf{B}}_{3} & \cdots & \overline{\mathbf{B}}_{N} \end{bmatrix} \begin{bmatrix} x_{I} \\ x_{2} \\ x_{3} \\ \vdots \\ x_{N} \end{bmatrix},$$
(18)

which can be expressed as:

$$\mathbf{Y} = \overline{\mathbf{C}}(\mathbf{M}\mathbf{X}),\tag{19}$$

where M represents the second term on the right-hand side232of Eq. 18. Recall from Eq. 3 that the result Y obtained by233self-attention processing can be expressed as:234

$$\mathbf{Y} = \mathbf{S}\mathbf{V} = (\mathbf{S}\boldsymbol{X})\mathbf{W}_V^\top \tag{20} \quad 235$$

From the perspective of self-attention, by comparing 236 Eqs. 19 and 20, the essence of Mamba is to generate a ma-237 trix M similar to similarity matrix S, such that the result 238 of $\mathbf{M}\mathbf{X}$ is based on the correlation between vectors of \mathbf{X} . 239 Although the final result of $\mathbf{M}\mathbf{X}$ is left multiplied by a map-240 ping matrix C, while the result of SX is right multiplied by 241 a mapping matrix \mathbf{W}_{V}^{\top} , the geometric meaning of the two 242 are the same. 243

3.2. Limitations of Mamba in video understanding 244

From the perspective of self-attention, the concept of245Mamba is similar: both use similarity matrices. We now246analyze the differences between the similarity matrices of247Mamba and self-attention, and discuss the limitations of248Mamba in the context of the video understanding task.249

221

229

230

¹The original SSM [13] employs h'(t) = Ah(t) + Bx(t), with h(t) the hidden state from previous time step t - 1, and h'(t) the updated current hidden state, replacing h(t). Considering this approach may lead to ambiguity, we have adopted the updated description.

309

310

Limitation 1: Historical decay. Matrix M in Eq. 19
corresponds to the second right-hand term in Eq. 18, which
is a lower triangular matrix of the form:

253 $\mathbf{M} = \begin{bmatrix} m_{11} & 0 & 0 & \cdots & 0 \\ m_{21} & m_{22} & 0 & \cdots & 0 \\ m_{31} & m_{32} & m_{33} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ m_{N1} & m_{N2} & m_{N3} & \cdots & m_{NN} \end{bmatrix}.$ (21)

By comparing \mathbf{M} with matrix \mathbf{S} in self-attention, we 254 255 find that outputs in Mamba favor more recent information, 256 because the more weights are zero, the earlier the token is observed. For example, for input $[x_1, x_2, x_3]$, the out-257 put Mx_1 in Mamba is $m_{11}x_1$ while the output Sx_1 is 258 $s_{11}x_1 + s_{12}x_2 + s_{13}x_3$ in self-attention. This indicates that, 259 260 in Mamba, the influence of earlier observed tokens on the final result is greatly diminished. We refer to this limitation 261 as historical decay. 262

In the NLP domain, more recent dialogue information 263 often has more impact on the final judgment, so this effect 264 265 is acceptable. However, in the computer vision domain, the order of the tokens has less meaning. Previous works such 266 267 as Vision Mamba [61] and VMamba [27] tried to solve this issue by processing the token sequence in both forward and 268 269 backward directions. This produces better results, but significant deficiencies still exist. We will explain this theoret-270 271 ically below.

272 When processing bi-directionally, the results generated 273 from input forward tokens $[x_1, \dots, x_N]$, denoted as $M_f X$, 274 and the results generated from input backward tokens 275 $[x_N, \dots, x_I]$, denoted as $M_b X$, are linearly combined to 276 generate the final result $M_{bi} X$ with M_{bi} a dense matrix. 277 As a result, the influence of historical information on the 278 result is increased, consequently leading to better results.

For example, for the input tokens $[x_1, x_2, x_3]$, $M_f X$ and $M_b X$ can be expressed as:

281
$$\mathbf{M}_{f} \boldsymbol{X} = \begin{bmatrix} f_{11} & 0 & 0 \\ f_{21} & f_{22} & 0 \\ f_{31} & f_{32} & f_{33} \end{bmatrix} \begin{bmatrix} \boldsymbol{x}_{I} \\ \boldsymbol{x}_{2} \\ \boldsymbol{x}_{3} \end{bmatrix} = \begin{bmatrix} h_{1f} \\ h_{2f} \\ h_{3f} \end{bmatrix}, \quad (22)$$

282

$$\mathbf{M}_{b}\boldsymbol{X} = \begin{bmatrix} b_{33} & 0 & 0 \\ b_{23} & b_{22} & 0 \\ b_{13} & b_{12} & b_{11} \end{bmatrix} \begin{bmatrix} \boldsymbol{x}_{3} \\ \boldsymbol{x}_{2} \\ \boldsymbol{x}_{1} \end{bmatrix} = \begin{bmatrix} h_{3b} \\ h_{2b} \\ h_{1b} \end{bmatrix}, \quad (23)$$

where f_{ij} represents the similarity score during the forward process, and b_{ij} is the similarity score in the backward direction. After bi-directional computation, with the outputs linearly combined, the results are expressed as:

$$h_{1} = h_{1f} + h_{1b} = f_{11}x_{I} + b_{13}x_{3} + b_{12}x_{2} + b_{11}x_{I}$$

$$h_{2} = h_{2f} + h_{2b} = f_{21}x_{I} + f_{22}x_{2} + b_{23}x_{3} + b_{22}x_{2} \quad (24)$$

$$h_{3} = h_{3f} + h_{3b} = f_{31}x_{I} + f_{32}x_{2} + f_{33}x_{3} + b_{33}x_{3}$$

We can write Eq. 24 in matrix form:

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \begin{bmatrix} f_{11} + b_{11} & b_{12} & b_{13} \\ f_{21} & f_{22} + b_{22} & b_{23} \\ f_{31} & f_{32} & f_{33} + b_{33} \end{bmatrix} \begin{bmatrix} \boldsymbol{x}_1 \\ \boldsymbol{x}_2 \\ \boldsymbol{x}_3 \end{bmatrix}$$

$$= \mathbf{M}_{bi} \begin{bmatrix} \boldsymbol{x}_1 \\ \boldsymbol{x}_2 \\ \boldsymbol{x}_3 \end{bmatrix}.$$

$$(25) \quad 289$$

The bi-directional computation transforms matrix M 290 from a lower triangular matrix to a dense matrix \mathbf{M}_{bi} , 291 thereby capturing more historical information and effectively avoiding historical decay. When extending to the case of N input tokens $[\mathbf{x}_1, \dots, \mathbf{x}_N]$, \mathbf{M}_{bi} can be written as: 294

$$\mathbf{M}_{bi} = \begin{bmatrix} f_{11} + b_{11} & b_{12} & b_{13} & \cdots & b_{1N} \\ f_{21} & f_{22} + b_{22} & b_{23} & \cdots & b_{2N} \\ f_{31} & f_{32} & f_{33} + b_{33} & \cdots & b_{3N} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ f_{N1} & f_{N2} & f_{N3} & \cdots & f_{NN} + b_{NN} \end{bmatrix}.$$
(26) 295

The diagonal elements of M_{bi} contain duplicates of the 296 similarity between a token and itself. For example, f_{33} 297 and b_{33} each represent the similarity between token x_3 and 298 itself. Consequently, the similarity is effectively doubled 299 which weakens the association with other tokens. One pos-300 sible approach is to adjust M_f and M_b using a weight 301 coefficient z through a linear combination. However, this 302 method is not very effective. For example, to prevent the 303 similarity of diagonal elements from being doubled, the z304 value is set to 0.5. However, after multiplying other ele-305 ments by 0.5, their weight values are significantly reduced, 306 which means that the association of each token with other 307 tokens is weakened. 308

Limitation 2: Element contradiction. By analyzing non-zero elements m_{ij} in M of Eq. 18, we derive:

$$m_{ij} = \mathbf{A}_i m_{i-1,j} \tag{27} \quad \mathbf{311}$$

After multiple iterations, the above equation results in im-312 plicit consideration of the correlation between previous to-313 kens and token j when computing the correlation between 314 token i and token j. As a result, m_{ij} exhibits stronger con-315 textual dependencies compared to the elements s_{ij} in the 316 matrix S. This might explain why Mamba achieves better 317 performance than transformers in the field of NLP. While 318 this is advantageous in the NLP domain, for the computer 319 vision domain, input tokens often lack semantic connec-320 tions. The consideration of the influence of other tokens 321 on each element can lead to significant drawbacks. 322

Interleaved token structures are common when processing images. Tokens that "belong together" might not be subsequently processed. For example, in an image classification task, input tokens $[x_1, x_2, x_3]$ might represent image regions [dog, other, dog]. Ideally, m_{31} should be high and 327



Figure 1. (a) Framework of VideoMambaPro with K bi-directional Mamba blocks. (b) In each bi-directional Mamba block, we employ forward residual SSM and masked backward Residual SSM.

 m_{21} low. Following Eq. 27, $m_{31} = \overline{\mathbf{A}}_3 m_{21}$, so $\overline{\mathbf{A}}_3$ needs to 328 329 increase. However, this causes $m_{32} = \overline{\mathbf{A}}_3 m_{22}$ to increase because m_{22} is also high. But, theoretically, m_{32} should 330 331 be low. This causes an element contradiction. Especially for video understanding, such contradictions are common 332 because most video regions contain background and other 333 334 irrelevant information, making relevant tokens sparse. Consequently, the performance of Mamba applied to video anal-335 336 ysis tasks is underwhelming [22, 24, 56].

337 4. VideoMambaPro

We use VideoMamba [22] as backbone, and propose two
adaptations to address historical decay and element contradiction. The resulting architecture is termed VideoMambaPro (VMP). With minor adjustments, our adaptations can
also be applied to related Mamba models.

To address historical decay, we keep the result of $M_f X$ but we use masked computation in the backward process. Specifically, we assign a mask to the diagonal elements of M_b , setting their values to 0, and then proceed with the calculations in Eqs 21–25. We thus eliminate the duplicate similarity on the diagonal, without affecting other elements. The final M_{bi} is expressed as:

350
$$M_{bi} = \begin{bmatrix} f_{11} & b_{12} & b_{13} & \cdots & b_{1N} \\ f_{21} & f_{22} & b_{23} & \cdots & b_{2N} \\ f_{31} & f_{32} & f_{33} & \cdots & b_{3N} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ f_{N1} & f_{N2} & f_{N3} & \cdots & f_{NN} \end{bmatrix} .$$
(28)

To solve element contradiction, we propose residual SSM. This solution is inspired by residual connections, to distribute the requirement for $\overline{\mathbf{A}}_i$ in m_{ij} over multiple $\overline{\mathbf{A}}_i$. This helps to avoid contradictions caused by interleaved sequence structures. For example, for our previous example input sequence $[x_1, x_2, x_3]$, which represents regions [dog, other, dog], we let $m_{31} = \overline{A}_3 m_{21} + \overline{A}_3$. This way, the requirement for a single \overline{A}_3 can be split into two parts, thus avoiding contradictions. This can be expressed as: 359

$$m_{ij} = \mathbf{A}_i m_{i-1,j} + \mathbf{A}_i \tag{29} \quad \mathbf{360}$$

We implement these solutions into VideoMamba [22], to 361 form VideoMambaPro (see Figure 1). Given input video 362 $X^{v} \in \mathbb{R}^{3 \times T \times H \times W}$, we first use a 3D convolution with a 363 $1 \times 16 \times 16$ kernel to convert \boldsymbol{X}^{v} into L non-overlapping 364 patch-wise tokens $X^p \in \mathbb{R}^{L \times C}$ with $L = t \times h \times w$ (t = t)365 $T, h = \frac{H}{16}, w = \frac{W}{16}$). Because SSM is sensitive to token 366 positions, and in line with [22], we include learnable spatial 367 and temporal position embeddings $p_s \in \mathbb{R}^{(hw+1) \times C}$ and 368 $p_t \in \mathbb{R}^{t \times C}$. Input tokens X are expressed as: 369

$$X = [X_{cls}, X] + p_s + p_t,$$
 (30) 370

376

where X_{cls} is a learnable classification token positioned at371the start of the sequence. Input tokens X pass through K372Mamba blocks, and the final layer's [CLS] token is used for373classification, after normalization and linear projection.374

5. Experiments 375

5.1. Experimental setup

Datasets.We evaluate VideoMambaPro on five video377benchmarks:(a) Kinetics-400 (K400, [3]) comprises378 ~ 240 K training and ~ 20 K validation videos, each with379an average duration of 10 seconds and categorized into380400 classes.(b) Something-Something V2 (SSv2, [10]) in-cludes ~ 160 K training and ~ 20 K validation videos with382

an average duration of 4 seconds, and 174 motion-centric classes. (c) UCF-101 [44] is a relatively small dataset, consisting of ~9.5K training and ~3.5K validation videos. (d) HMDB51 [18] is also a compact video dataset, containing ~3.5K training and ~1.5K validation videos. (e) AVA [14] is a dataset for spatio-temporal localization of human actions with ~211k and ~57k validation video segments.

390 Implementation. In line with VideoMamba, we introduce 391 three models with increasing embedding dimension and number of bi-directional Mamba blocks K: Tiny, Small, 392 and Middle (details in supplementary material). To com-393 pare with VideoMamba, we pre-train VideoMambaPro on 394 395 ImageNet-1K (IN-1K). On K400, we also pre-train with IN-1K, fine-tune on the training set and report on the validation 396 set. For K400, we also report on the larger 336^2 input size. 397 During pre-training, we follow DeiT [48] by applying a cen-398 ter crop to obtain the 224^2 sized images. We apply random 399 cropping, random horizontal flipping, label-smoothing reg-400 401 ularization, mix-up, and random erasing as data augmentations. We use AdamW [28] with a momentum of 0.9, a 402 batch size of 1024, and a weight decay of 0.05. We employ 403 a cosine learning rate schedule during training, 1×10^{-3} 404 initial learning rate over 300 epochs. The fine-tuning set-405 406 tings follow VideoMAE [47]. We resize frames to 224^2 , 407 and use AdamW with a momentum of 0.9 and a batch size of 512. The evaluation process is the same as VideoMamba, 408 and more details are in the supplementary materials. These 409 materials also include results on IN-1K image classification. 410

Method	Pre-train	Input	Crops	Param	FLOP	Top1	Top5
MViTv1-B [6]		32×224^2	5×1	37M	350G	80.2	94.4
MViTv2-S [25]		16×224^2	5×1	35M	320G	81.0	94.6
Uniformer-S [20]	IN-1K	16×224^2	4×1	21M	168G	80.8	94.7
Uniformer-B [20]	IN-1K	16×224^2	4×1	50M	388G	82.0	95.1
Uniformer-B [20]	IN-1K	32×224^2	4×3	50M	3.1T	83.0	95.4
STAM [42]	IN-21K	64×224^2	1×1	121M	1.0T	79.2	-
TimeSformer-L [2]	IN-21K	96×224^2	1×3	121M	7.1T	80.7	94.7
ViViT-L [1]	IN-21K	16×224^2	4×3	311M	47.9T	81.3	94.7
Mformer-HR [36]	IN-21K	16×336^2	10×3	311M	28.8T	81.1	95.2
VideoMAE-H [47]	IN-21K	16×224^2	5×3	633M	17.9T	86.6	97.1
X-CLIP-L/14 [32]	CLIP-400M	16×336^2	4×3	453M	37.0T	87.7	_
MTV-H [55]	$60M^{1}$	32×224^2	4×3	1120M	44.5T	89.1	98.2
InternVideo-1B [51]	$412M^{2}$	64×224^2	16×4	1300M	86.2T	91.1	98.9
InternVideo2-1B [52]	$414M^{3}$	16×224^2	16×4	1000M	_	91.6	_
InternVideo2-6B [52]	$414M^3$	16×224^2	16×4	5903M	_	92.1	
VideoMamba-Ti	IN-1K	32×224^2	4×3	7M	0.4T	78.8	93.9
VideoMamba-Ti	IN-1K	64×384^2	4×3	7M	2.4T	80.3	94.8
VideoMamba-S	IN-1K	32×224^2	4×3	26M	1.6T	81.5	95.2
VideoMamba-S	IN-1K	64×384^2	4×3	26M	4.7T	82.7	95.6
VideoMamba-M	IN-1K	32×224^2	4×3	74M	4.8T	82.4	95.7
VideoMamba-M	IN-1K	64×384^2	4×3	74M	28.4T	83.3	96.1
VideoMambaPro-Ti	IN-1K	32×224^2	4×3	7M	0.4T	81.6	95.9
VideoMambaPro-Ti	IN-1K	64×384^2	4×3	7M	2.2T	83.3	96.1
VideoMambaPro-S	IN-1K	32×224^2	4×3	25M	1.6T	83.3	96.0
VideoMambaPro-S	IN-1K	64×384^2	4×3	25M	4.4T	84.5	96.6
VideoMambaPro-M	IN-1K	32×224^2	4×3	72M	4.7T	84.0	96.4
VideoMambaPro-M	IN-1K	64×384^2	4×3	72M	27.0T	85.0	96.7

Table 1. Performance on K400. Top part of the table are Transformer models, bottom part are Mamba models. We report crops (temporal \times spatial) and FLOPs for inference. —: not reported. ¹ IN-21K+WTS

5.2. Comparison with state-of-the-art

K400. Results appear in Table 1. Compared to Video-412 Mamba, VideoMambaPro has slightly fewer parameters and 413 FLOPs. This is primarily because VideoMamba employs an 414 additional projection layer to generate the weight coefficient 415 z to adjust \mathbf{A}_f and \mathbf{A}_b . See the supplementary materials for 416 an architecture comparison. VideoMambaPro outperforms 417 VideoMamba across model and input sizes. With 224² in-418 puts and pre-trained only on IN-1K, the best-performing 419 VideoMambaPro-M achieves a top-1 accuracy of 84.0%, 420 1.6% higher than VideoMamba-M. Further comparisons ap-421 pear in Section 5.4. Increasing the input size to 336^2 leads 422 to a performance improvement of 1.0-1.7%. 423

VideoMambaPro scores lower than the recent 424 InternVideo2-1B [52] by 7.6%, but was only pre-trained on 425 IN-1K and has significantly fewer parameters (1000M vs 426 72M) and inference only takes ~5.5% of the FLOPs. 427

Method	Pre-train	Input	Crops	Param	FLOP	Top1	Top5
MViTv1-B [6]	K400	16×224^2	1×3	37M	213G	64.7	89.2
MViTv1-B [6]	K400	32×224^2	1×3	37M	510G	67.1	90.8
MViTv2-S [25]	K400	16×224^2	1×3	35M	195G	68.2	91.4
MViTv2-B [25]	K400	32×224^2	1×3	51M	675G	70.5	92.7
Uniformer-S [20]	IN-1K+K400	16×224^2	1×3	21M	126G	67.7	91.4
Uniformer-B [20]	IN-1K+K400	16×224^2	1×3	50M	291G	70.4	92.8
TimeSformer-L [2]	IN-21K	16×224^2	1×3	121M	5.1T	62.5	-
ViViT-L [1]	IN-21K+K400	$16 imes 224^2$	4×3	311M	47.9T	65.4	89.8
Mformer-HR [36]	IN-21K+K400	16×336^2	1×3	311M	3.6T	68.1	91.2
MaskFeat-L [53]	IN-21K	$64 imes 312^2$	4×3	218M	8.5T	75.0	95.0
VideoMAE-L [47]	IN-21K	32×224^2	1×3	305M	4.3T	75.4	95.2
TubeViT-L [39]	IN-1K	32×224^2	4×3	311M	9.5T	76.1	95.2
InternVideo-1B [51]	See Table 1	64×224^2	16×4	1300M	86.2T	77.2	95.9
InternVideo2-1B [52]	See Table 1	64×224^2	16×4	1000M		77.1	_
InternVideo2-6B [52]	See Table 1	64×224^2	16×4	5903M	-	77.4	-
VideoMamba-Ti	IN-1K	16×224^2	2×3	7M	102G	66.0	89.6
VideoMamba-S	IN-1K	16×224^2	2×3	26M	408G	67.6	90.9
VideoMamba-M	IN-1K	16×224^2	4×3	74M	2.4T	68.3	91.4
VideoMambaPro-Ti	IN-1K	$16 imes 224^2$	2×3	7M	96G	67.9	91.2
VideoMambaPro-S	IN-1K	16×224^2	2×3	25M	382G	68.8	91.4
VideoMambaPro-M	IN-1K	16×224^2	4×3	72M	2.2T	69.4	91.6

Table 2. Performance on SSv2. —: not reported. Top part of the table are Transformer models, bottom part are Mamba models.

SSv2. Results appear in Table 2. VideoMambaPro 428 outperforms VideoMamba by 1.1-1.9%. It also out-429 performs several popular transformer models. Although 430 InternVideo-1B [51] and InternVideo2-6B [52] outperform 431 our VideoMambaPro-M by 7.8% and 8.0%, respectively, 432 they require 18.0-82 times more parameters and at least 39 433 times more FLOPs. Again, we expect that the performance 434 for VideoMambaPro will increase with more pre-training. 435

UCF-101/HMDB51/AVA V2.2. From Table 3, it shows
that VideoMambaPro-M is competitive, and outperforms436VideoMamba by 3.4% and 1.8% on UCF-101 and
HMDB51, respectively. VideoMambaPro-M achieves 31.9
mAP on AVA V2.2, which is 10.7% lower than VideoMAE
V2 [50] but with an order of magnitude fewer parameters
and FLOPs and pre-trained only on IN-1K (see Table 4).436

² CLIP-400M+WebVid+HowTo+K710+SSv2+AVA2.2+more.

³ LAION-300M+KMash+WebVid+InternVid+LLaVA+more.

Method	Params	UCF-101	HMDB51
VideoMoCo [33]	15M	78.7	49.2
CoCLR [16]	9M	81.4	52.1
MemDPC [15]	32M	86.1	54.5
Vi ² CLR [4]	9M	89.1	55.7
VideoMAE [47]	87M	91.3	62.6
GDT [35]	33M	95.2	72.8
VideoMAE V2 [50]	1050M	99.6	88.1
VideoMamba-M	74M	88.2	60.8
VideoMambaPro-M	72M	91.6	63.2

Table 3. Results on UCF-101 and HMDB51.

Method	FLOPs	Param	mAP
SlowFast R101 [8]	138G	53M	23.8
VideoMAE-B [47]	180G	87M	26.7
MViTv2-B [25]	225G	51M	30.5
ObjectTransformer [54]	243G	86M	31.0
MViTv2-L [25]	2828G	213M	34.4
ST-MAE-H [9]	1193G	632M	36.2
VideoMAE V2 [50]	4220G	1050M	42.6
VideoMamba-M [19]	202G	74M	30.1
VideoMambaPro-M	183G	72M	31.9

Table 4. Results on AVA V2.2.

Models	Input	Top-1	Top-5
VideoMamba-M (baseline)	32×224^2	82.4	95.7
VideoMambaPro-M (w/o residual)	32×224^2	83.6 (+1.2)	96.0 (+0.3)
VideoMambaPro-M (w/o masking)	32×224^2	83.0 (+1.0)	95.8 (+0.1)
VideoMambaPro-M	32×224^2	84.0 (+1.6)	96.4 (+0.7)

Table 5. Ablation study on K400, with and without masked backward computation and elemental residual connections.

443 5.3. Ablation study

Influence of masked backward computation and ele-444 mental residual connection. We have identified two lim-445 446 itations that exist in VideoMamba: historical decay and element contradiction. We introduced masked backward 447 448 computation and elemental residual connections to address these respective issue. Here, we analyze the impact of 449 450 each solution. We use the same settings as before, with VideoMambaPro-M and pre-training on IN-1K. We sum-451 marize the performance of VideoMambaPro-M on K400 in 452 453 Table 5. Both solutions contribute to an improved score, and their effect is partly complementary. This indicates that 454 the two limitations exist simultaneously in VideoMamba. 455

Evidence of limitations of Mamba. We have theo-456 457 retically identifies historical decay and element contradiction as fundamental problems in Mamba, and propose the 458 masked backward computation and elemental residual con-459 nection to solve the limitations. Previous experiments on 460 image and video action recognition demonstrated the ef-461 462 fectiveness of our approach, and we further provide con-463 crete evidence that the increased performance stems from addressing these two issues.

We performed two tests on ImageNet-1K image clas-465 sification task to build the connection between the issues 466 in practice and our approach. First, to investigate whether 467 the residual connection alleviates element contradiction, we 468 randomly replaced a percentage of the patches by randomly 469 selected patches from other classes, which will increase el-470 ement contradiction because more irrelevant information is 471 present. The table 6 shows that the relative performance 472 of VideoMambaPro over VideoMamba increases to 9.2% 473 when 20% of the patches is replaced. For higher percent-474 ages, the gap is reduced along with the overall score. This 475 demonstrates that irrelevant (or even misleading) informa-476 tion can better be dealt with using our approach.

		Replacement ratio					
	0	5	10	20	30	50	80
VideoMamba-Ti	76.9	73.2	68.0	59.4	52.9	47.6	28.9
VideoMambaPro-Ti	78.9	78.0	76.1	68.6	55.9	48.2	29.1

Table 6. Top-1 accuracy of VideoMambaPro-Ti and VideoMamba-Ti with various replacement ratios (in %) on ImageNet-1K.

Second, to verify that masking addresses the issue of historical decay, we adopted a progressive masking approach, gradually decreasing the mask values of the diagonal elements of M_b from 1 to 0, and then retrained the model 481 on ImageNet-1K. The results below show that VideoMambaPro's accuracy gradually increases when more masking 483 is applied, evidencing the merits of the approach. 484

		Mask	value	
	1	0.7	0.5	0
VideoMambaPro-Ti	76.9	77.0	77.2	77.8

Table 7. Top-1 accuracy across mask values on ImageNet-1K.

5.4. Comparison with VideoMamba on K400

We more thoroughly compare the differences between 486 VideoMamba and VideoMambaPro by investigating the rel-487 ative performance per class. We then present a statistical 488 comparison between the results of both backbones. 489 490

Class analysis. We compare VideoMambaPro-M with 224^2 image size pre-trained on IN-1K to a VideoMamba-M baseline with the same settings. We show the relative performance for all classes of K400 in Figure 2. For over 95% of the classes, VideoMambaPro shows improvement. Although there is a lower performance for certain classes, the decrease is typically limited.

The majority of the classes sees a $\sim 1.8\%$ improvement, 497 which is substantial. For a small number of classes, Video-498 MambaPro performs >2% better than VideoMamba. Only 499 a fraction of the classes is negatively affected by the solu-500 tions introduced in VideoMambaPro. 501

485

491

492

493

494

495



Figure 2. Relative accuracy per class on Kinetics-400 by comparing VideoMambaPro-M to VideoMamba-M.

508

Statistical comparison. In order to understand whether the improvements of VideoMambaPro over VideoMamba are statistically significant, we compare the K400 results of the respective Middle models, both pre-trained on IN-1K and size 224×224 . Other settings are also the same. For each test sample, we check whether if it correctly classified by either model. The results appear in Table 8.

		VideoMambaPro-M				
		True	False	Total		
	True	14,302	469	14,771		
VideoMamba-M	False	1,833	3,302	5,135		
	Total	16,135	3,771	19,906		

Table 8. Contingency table for K400 test items for VideoMamba-M and VideoMambaPro-M.

We used the McNemar test, a non-parametric test with 509 510 a single degree of freedom. It checks whether the number of items that are incorrectly classified by VideoMam-511 baPro (n_{10}) but not VideoMamba is substantially lower 512 than the number of items misclassified by VideoMamba 513 514 but not VideoMambaPro (n_{01}) . The test is calculated as $\chi^2 = \frac{(n_{01} - n_{10})^2}{(n_{01} + n_{10})}$. The resulting value of 808.2 corresponds 515 to a significance level of p < 0.001. We can thus conclude 516 that VideoMambaPro-M is statistically significantly better 517 than VideoMamba-M. Because we relied on the aggregated 518 519 performance reported in papers for other methods, we cannot report statistical comparisons here. 520

521 5.5. Computation cost analysis

Finally, we compare the performance of VideoMambaPro
with various model sizes with other approaches on K400.
We map the top-1 to the number of parameters and FLOPs
in Figures 3 and 4, respectively. VideoMambaPro performs favorably in terms of the accuracy-compute trade-



Figure 3. Top-1 accuracy versus number of parameters of Video-MambaPro and other models on Kinetics-400.



Figure 4. Top-1 accuracy versus number of FLOPs of VideoMambaPro and other models on Kinetics-400.

off. Importantly, VideoMambaPro was trained on much less527data than other models, which might provide opportunities528for further accuracy gains without additional compute and529memory requirements.530

6. Conclusion

From a mathematical comparison with self-attention, we 532 have identified two limitations in how Mamba processes to-533 ken sequences. We argue that these limitations constrain 534 Mamba's potential, especially in video understanding tasks. 535 To address the two limitations, we have introduced Video-536 MambaPro (VMP). It takes VideoMamba and introduces 537 the masked backward State Space Model (SSM), and adds 538 residual connections in both forward and backward SSM. 539 In experiments on Kinetics-400, Something-Something V2, 540 HMDB51, UCF-101, and AVA V2.2, VideoMambaPro 541 consistently demonstrates improved performance over the 542 vanilla VideoMamba. We expect that extensive pre-training 543 further elevates the performance of Mamba models for 544 video tasks, making it an increasingly attractive, efficient 545 alternative to large transformer models. 546

558

559

560

580

581

582

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

547 References

- 548 [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen
 549 Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vi550 sion transformer. In *ICCV*, pages 6836–6846, 2021. 6
- [2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is
 space-time attention all you need for video understanding?
 In *ICML*, pages 813–824, 2021. 6
- [3] Joao Carreira and Andrew Zisserman. Quo vadis, action
 recognition? A new model and the Kinetics dataset. In
 CVPR, pages 6299–6308, 2017. 5
 - [4] Ali Diba, Vivek Sharma, Reza Safdari, Dariush Lotfi, Saquib Sarfraz, Rainer Stiefelhagen, and Luc Van Gool. Vi2CLR: Video and image for visual contrastive learning of representation. In *CVPR*, pages 1502–1512, 2021. 7
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov,
 Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is
 worth 16x16 words: Transformers for image recognition at
 scale. In *ICLR*, 2021. 1
- [6] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li,
 Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*,
 pages 6824–6835, 2021. 6
- [7] Yuxin Fang, Bencheng Liao, Xinggang Wang, Jiemin Fang,
 Jiyang Qi, Rui Wu, Jianwei Niu, and Wenyu Liu. You
 only look at one sequence: Rethinking transformer in vision
 through object detection. *NeurIPS*, 34:26183–26197, 2021.
 2
- 577 [8] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and
 578 Kaiming He. Slowfast networks for video recognition. In
 579 *ICCV*, pages 6202–6211, 2019. 7
 - [9] Christoph Feichtenhofer, Haoqi Fan, Yanghao Li, and Kaiming He. Masked autoencoders as spatiotemporal learners. *NeurIPS*, 35:35946–35958, 2022. 2, 7
- [10] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim,
 Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz
 Mueller-Freitag, et al. The "Something Something" video
 database for learning and evaluating visual common sense.
 In *ICCV*, pages 5842–5850, 2017. 5
- [11] Albert Gu and Tri Dao. Mamba: Linear-time sequence
 modeling with selective state spaces. arXiv preprint
 arXiv:2312.00752, 2023. 1, 2
- 592 [12] Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christo 593 pher Ré. Hippo: Recurrent memory with optimal polynomial
 594 projections. *NeurIPS*, 33:1474–1487, 2020. 2
- [13] Albert Gu, Karan Goel, and Christopher Ré. Efficiently
 modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021. 1, 2, 3
- [14] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, et al.
 AVA: A video dataset of spatio-temporally localized atomic visual actions. In *CVPR*, pages 6047–6056, 2018. 6

- [15] Tengda Han, Weidi Xie, and Andrew Zisserman. Memoryaugmented dense predictive coding for video representation learning. In *ECCV*, pages 312–329, 2020. 7
 603
 604
 605
- [16] Tengda Han, Weidi Xie, and Andrew Zisserman. Selfsupervised co-training for video representation learning. Advances in neural information processing systems, 33:5679– 5690, 2020. 7
- [17] Feyza Duman Keles, Pruthuvi Mahesakya Wijewardena, and Chinmay Hegde. On the computational complexity of selfattention. In *International Conference on Algorithmic Learning Theory*, pages 597–619. PMLR, 2023. 1
- [18] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. HMDB: a large video database for human motion recognition. In *ICCV*, pages 2556–2563, 2011. 6
- [19] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Limin Wang, and Yu Qiao. UniformerV2: Spatiotemporal learning by arming image vits with video uniformer. arXiv preprint arXiv:2211.09552, 2022. 1, 7
- [20] Kunchang Li, Yali Wang, Gao Peng, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. Uniformer: Unified transformer for efficient spatial-temporal representation learning. In *ICLR*, 2022. 6
- [21] Kunchang Li, Yali Wang, Yizhuo Li, Yi Wang, Yinan He, Limin Wang, and Yu Qiao. Unmasked teacher: Towards training-efficient video foundation models. In *ICCV*, pages 19948–19960, 2023. 2
- [22] Kunchang Li, Xinhao Li, Yi Wang, Yinan He, Yali Wang, Limin Wang, and Yu Qiao. VideoMamba: State space model for efficient video understanding. arXiv preprint arXiv:2403.06977, 2024. 1, 2, 5
- [23] Shufan Li, Harkanwar Singh, and Aditya Grover. Mamband: Selective state space modeling for multi-dimensional data. In *European Conference on Computer Vision*, pages 75–92. Springer, 2024. 2
- [24] Wenrui Li, Xiaopeng Hong, and Xiaopeng Fan. Spikemba: Multi-modal spiking saliency mamba for temporal video grounding. arXiv preprint arXiv:2404.01174, 2024. 5
- [25] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. MViTv2: Improved multiscale vision transformers for classification and detection. In *CVPR*, pages 4804–4814, 2022. 6, 7
- [26] Qin Liu, Zhenlin Xu, Gedas Bertasius, and Marc Niethammer. SimpleClick: Interactive image segmentation with simple vision transformers. In *ICCV*, pages 22290–22300, 2023.
 2
- [27] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and Yunfan Liu. VMamba: Visual state space model. *arXiv preprint arXiv:2401.10166*, 2024. 2, 4
- [28] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [29] Hui Lu, Hu Jian, Ronald Poppe, and Albert Ali Salah. Enhancing video transformers for action understanding with vlm-aided training. *arXiv preprint arXiv:2403.16128*, 2024.
 1

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

- [30] Jiachen Lu, Jinghan Yao, Junge Zhang, Xiatian Zhu, Hang
 Xu, Weiguo Gao, Chunjing Xu, Tao Xiang, and Li Zhang.
 Soft: Softmax-free transformer with linear complexity. *NeurIPS*, 34:21297–21309, 2021. 2
- [31] Harsh Mehta, Ankit Gupta, Ashok Cutkosky, and Behnam
 Neyshabur. Long range language modeling via gated state
 spaces. *arXiv preprint arXiv:2206.13947*, 2022. 1
- [32] Bolin Ni, Houwen Peng, Minghao Chen, Songyang Zhang,
 Gaofeng Meng, Jianlong Fu, Shiming Xiang, and Haibin
 Ling. Expanding language-image pretrained models for general video recognition. In *ECCV*, pages 1–18, 2022. 6
- [33] Tian Pan, Yibing Song, Tianyu Yang, Wenhao Jiang, and Wei
 Liu. VideoMoCo: Contrastive video representation learning with temporally adversarial examples. In *CVPR*, pages
 11205–11214, 2021. 7
- [34] Jinyoung Park, Hee-Seon Kim, Kangwook Ko, Minbeom
 Kim, and Changick Kim. Videomamba: Spatio-temporal selective state space model. In *European Conference on Com- puter Vision*, pages 1–18. Springer, 2024. 2
- 679 [35] Mandela Patrick, Yuki Markus Asano, Ruth Fong, João F.
 680 Henriques, Geoffrey Zweig, and Andrea Vedaldi. Multi681 modal self-supervision from generalized data transforma682 tions. *arXiv preprint arXiv:2003.04298*, 2020. 7
- [36] Mandela Patrick, Dylan Campbell, Yuki Asano, Ishan Misra,
 Florian Metze, Christoph Feichtenhofer, Andrea Vedaldi,
 and Joao F Henriques. Keeping your eye on the ball: Trajectory attention in video transformers. *Advances in neural information processing systems*, 34:12493–12506, 2021. 6
- [37] Badri Narayana Patro and Vijay Srinivas Agneeswaran.
 Mamba-360: Survey of state space models as transformer
 alternative for long sequence modelling: Methods, applications, and challenges. *arXiv preprint arXiv:2404.16112*,
 2024. 1
- [38] Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak,
 Samuel Arcadinho, Huanqi Cao, Xin Cheng, Michael
 Chung, Matteo Grella, Kranthi Kiran GV, et al. RWKV:
 Reinventing RNNs for the transformer era. *arXiv preprint arXiv:2305.13048*, 2023. 1, 2
- [39] AJ Piergiovanni, Weicheng Kuo, and Anelia Angelova. Rethinking video ViTs: Sparse video tubes for joint image and video learning. In *CVPR*, pages 2214–2224, 2023. 6
- [40] Hongyu Ren, Hanjun Dai, Zihang Dai, Mengjiao Yang, Jure
 Leskovec, Dale Schuurmans, and Bo Dai. Combiner: Full attention transformer with sparse computation cost. *NeurIPS*, 34:22470–22482, 2021. 2
- [41] Chaitanya Ryali, Yuan-Ting Hu, Daniel Bolya, Chen Wei, Haoqi Fan, Po-Yao Huang, Vaibhav Aggarwal, Arkabandhu Chowdhury, Omid Poursaeed, Judy Hoffman, et al. Hiera: A hierarchical vision transformer without the bells-andwhistles. *arXiv preprint arXiv:2306.00989*, 2023. 1
- [42] Gilad Sharir, Asaf Noy, and Lihi Zelnik-Manor. An image is
 worth 16x16 words, what is a video worth? *arXiv preprint arXiv:2103.13915*, 2021. 6
- [43] Yuanyuan Shen, Edmund M-K Lai, and Mahsa Mohaghegh.
 Effects of similarity score functions in attention mechanisms on the performance of neural question answering systems. *Neural Processing Letters*, 54(3):2283–2302, 2022. 2

- [44] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah.
 UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
 6
- [45] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *ICCV*, pages 7262–7272, 2021. 2
- [46] Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. Retentive network: A successor to transformer for large language models. arXiv preprint arXiv:2307.08621, 2023. 1, 2
- [47] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. VideoMAE: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *NeurIPS*, 35:10078– 10093, 2022. 1, 2, 6, 7
- [48] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, pages 10347–10357. PMLR, 2021. 6
- [49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017. 1, 2
- [50] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao. VideoMAE V2: Scaling video masked autoencoders with dual masking. In *CVPR*, pages 14549–14560, 2023. 1, 6, 7
- [51] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, et al. InternVideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022. 1, 6
- [52] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Jilan Xu, Zun Wang, et al. InternVideo2: Scaling video foundation models for multimodal video understanding. arXiv preprint arXiv:2403.15377, 2024. 6
- [53] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked feature prediction for self-supervised visual pre-training. In *CVPR*, pages 14668–14678, 2022. 2, 6
- [54] Chao-Yuan Wu and Philipp Krahenbuhl. Towards long-form video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1884–1894, 2021. 7
- [55] Shen Yan, Xuehan Xiong, Anurag Arnab, Zhichao Lu, Mi Zhang, Chen Sun, and Cordelia Schmid. Multiview transformers for video recognition. In *CVPR*, pages 3333–3343, 2022. 6
- [56] Hanwei Zhang, Ying Zhu, Dan Wang, Lijun Zhang, Tianxiang Chen, and Zi Ye. A survey on Visual Mamba. arXiv preprint arXiv:2404.15956, 2024. 5
- [57] Zheng Zhang and Kil To Chong. Comparison between firstorder hold with zero-order hold in discretization of inputdelay nonlinear systems. In 2007 International Conference on Control, Automation and Systems, pages 2892–2896, 2007. 3

- [58] Zixiao Zhang, Xiaoqiang Lu, Guojin Cao, Yuting Yang,
 Licheng Jiao, and Fang Liu. ViT-YOLO: Transformer-based
 YOLO for object detection. In *ICCV*, pages 2799–2808,
 2021. 2
- [59] Chen Zhu, Wei Ping, Chaowei Xiao, Mohammad Shoeybi,
 Tom Goldstein, Anima Anandkumar, and Bryan Catanzaro.
 Long-short transformer: Efficient transformers for language
 and vision. *NeurIPS*, 34:17723–17736, 2021. 2
- [60] Lei Zhu, Xinjiang Wang, Zhanghan Ke, Wayne Zhang, and
 Rynson WH Lau. BiFormer: Vision transformer with bilevel routing attention. In *CVPR*, pages 10323–10333, 2023.
 2
- [61] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang,
 Wenyu Liu, and Xinggang Wang. Vision Mamba: Efficient
 visual representation learning with bidirectional state space
 model. *arXiv preprint arXiv:2401.09417*, 2024. 1, 2, 4

Snakes and Ladders: Two Steps Up for VideoMamba Supplementary Material

001We provide the architectures for VideoMambaPro mod-002els in Section 1. A comparison between the architectures003of VideoMamba and VideoMambaPro appears in Section 2.004Training details are presented in Section 3. Finally, we re-005port ImageNet-1K image classification results in Section 4,006to illustrate the potental of the proposed solutions for image007classification tasks.

1. VideoMambaPro architectures

We present the architecture details of VideoMambaProTiny (Ti), -Small (S), and -Middle (M) in Tables 1–3. The
differences are in the embedding dimension (192, 384, 576)
and the number of SSM blocks (24, 24, 32).

Stage	Tiny				
Patch Embedding	nn.Conv3d (kernel size = $16 \times 16 \times 1$, embedding dimension = 192)				
SSM	$\begin{bmatrix} MLP(768) \\ MLP(3072) \\ MHA (head = 12) \end{bmatrix} \times 24$				
Projection	Layer Normalization Dropout (ratio) Linear layer (1000) Softmax				

Table 1. Architecture details of VideoMambaPro-Ti.

Stage	Small				
Patch Embedding	nn.Conv3d (kernel size = $16 \times 16 \times 1$, embedding dimension = 384)				
SSM	$\begin{bmatrix} MLP(768) \\ MLP(3072) \\ MHA (head = 12) \end{bmatrix} \times 24$				
Projection	Layer Normalization Dropout (ratio) Linear layer (1000) Softmax				

Table 2. Architecture details of VideoMambaPro-S.

013 2. Architecture comparison with VideoMamba

We compare the architectures of VideoMambaPro andVideoMamba [1] in Figure 1. VideoMambaPro does not

Stage	Middle				
Patch Embedding	nn.Conv3d (kernel size = $16 \times 16 \times 1$, embedding dimension = 576)				
SSM	$\begin{bmatrix} MLP(768) \\ MLP(3072) \\ MHA (head = 12) \end{bmatrix} \times 32$				
Projection	Layer Normalization Dropout (ratio) Linear layer (1000) Softmax				

Table 3. Architecture details of VideoMambaPro-M.

have the linear layer to generate parameters z. Additionally,016our residual SSM and mask scheme do not introduce addi-
tional parameters or computational overhead, so our method017018019



Figure 1. Comparison between the bi-directional VideoMamba (top) and VideoMambaPro (bottom) blocks.

3. Implementation details

We conduct pre-training on ImageNet-1K and fine-tuning021on the Something-Something V2 and Kinetics-400 datasets022with 16 NVIDIA A100-80G GPUs.Models for UCF101023

1

and HMDB51 are trained with 8 A100-80G GPUs. The ex-024 periments on AVA V2.2 are conducted with 32 A100-80G 025 GPUs. The values of the hyperparameters are largely sim-026 ilar to those used in VideoMamba [1]. We linearly scale 027 028 the base learning rate with respect to the overall batch size, $lr = lr_{base} \times batchsize/256$. The pre-training details are 029 shown in Table 4, and the fine-tuning details on the other 030 datasets are listed in Tables 5-8. 031

config	image size: 224×224
optimizer	AdamW
base learning rate	1.5e-4
weight decay	0.1 (Tiny), 0.05 (Small, Middle)
minimal learning rate	1.0e-6
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.95$
batch size	512
learning rate schedule	cosine decay
warmup epochs	5 (Tiny), 10 (Small), 40 (Middle)
dropout ratio	0 (Tiny), 0.15 (Small), 0.5 (Middle)
augmentation	MultiScaleCrop
label smoothing	0.1

Table 4. Pre-training setting on ImageNet-1K

config	image size: 224×224			
optimizer	AdamW			
base learning rate	1.5e-4			
weight decay	0.1 (Tiny), 0.05 (Small, Middle)			
minimal learning rate	1.0e-6			
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.99$			
batch size	256			
learning rate schedule	cosine decay			
warmup epochs	5 (Tiny), 5 (Small) 10 (Middle)			
dropout ratio	0.1 (Tiny), 0.35 (Small), 0.6 (Middle)			
augmentation	RandAug (7, 0.25) (Tiny), RandAug (9, 0.5) (Small, Middle)			
label smoothing	0.1			
flip augmentation	yes			

Table 5. Fine-tuning setting for Kinetics-400

config	image size: 224×224			
optimizer	AdamW			
base learning rate	4e-4			
weight decay	0.1 (Tiny), 0.05 (Small, Middle)			
minimal learning rate	1.0e-6			
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.999$			
batch size	256			
learning rate schedule	cosine decay			
warmup epochs	5 (Tiny), 5 (Small) 10 (Middle)			
dropout ratio	0.1 (Tiny), 0.35 (Small), 0.6 (Middle)			
augmentation	RandAug (7, 0.25) (Tiny), RandAug (9, 0.5) (Small, Middle)			
label smoothing	0.1			
flip augmentation	no			

Table 6. Fine-tuning setting for Something-Something V2

4. Results on ImageNet-1K

We argue that the proposed solutions in handling the fea ture extraction capabilities of Mamba models are most ef fective when relevant tokens are more sparsely distributed

config	image size: 224×224			
optimizer	AdamW			
base learning rate	4e-4			
weight decay	0.1 (Tiny), 0.05 (Small, Middle)			
minimal learning rate	1.0e-6			
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.99$			
batch size	128			
learning rate schedule	cosine decay			
warmup epochs	5 (Tiny), 5 (Small) 10 (Middle)			
dropout ratio	0.1 (Tiny), 0.35 (Small), 0.6 (Middle)			
augmentation	RandAug (7, 0.25) (Tiny), RandAug (9, 0.5) (Small, Middle)			
label smoothing	0.1			
flip augmentation	yes			

Table 7. Fine-tuning setting for UCF101/HMDB51

config	image size: 224×224				
optimizer	AdamW				
base learning rate	1.5e-3 (Tiny), 2.5e-4 (Small, Middle)				
weight decay	0.051 (Tiny, Small, Middle)				
minimal learning rate	1.0e-6				
optimizer momentum	$\beta_1, \beta_2 = 0.9, 0.999$				
batch size	128				
learning rate schedule	cosine decay				
warmup epochs	5 (Tiny), 5 (Small) 10 (Middle)				
dropout ratio	0.1 (Tiny), 0.35 (Small) 0.6 (Middle)				
augmentation	RandAug (7, 0.25) (Tiny), RandAug (9, 0.5) (Small, Middle)				
label smoothing	0.1				
flip augmentation	yes				



in the input. While our main focus is on the video do-036 main, we also summary our experiments on image classi-037 fication. We pre-train VideoMambaPro on ImageNet-1K, 038 which contains 1.28M training images and 50K validation 039 images across 1,000 categories. All models are trained on 040 the training set, and top-1 accuracy on the validation set is 041 reported. For fair comparison, we adopt the same method 042 as VideoMamba, and our training settings primarily follows 043 DeiT [2]. When training on 224^2 input images, we use 044 AdamW with a momentum of 0.9 and a total batch size of 045 512. Training is performed on 8 A800 GPUs, with more 046 details provided in Table 4. The results are summarized in 047 Table 9. VideoMambaPro achieves accuracy gains of 0.9-048 2.0% over VideoMamba. 049

	_	_		
	Input	Param	FLOPs	Top-1
VideoMamba (Ti)	224 ²	7M	1.1G	76.9
VideoMambaPro (Ti)	$224^{\ 2}$	7M	1.1G	78.9
VideoMamba (S)	224 ²	26M	4.3G	81.2
VideoMambaPro (S)	224^{2}	25M	4.2G	82.4
VideoMamba (M)	224 ²	74M	12.7G	82.8
VideoMambaPro (M)	224^{2}	72M	12.4G	83.8
VideoMamba (M)	$448^{\ 2}$	75M	50.4G	83.8
VideoMambaPro (M)	$448^{\ 2}$	73M	49.6G	84.7

Table 9. ImageNet-1K pre-training results for VideoMamba and VideoMambaPro.

ICCV

#8505

050 References

- [1] Kunchang Li, Xinhao Li, Yi Wang, Yinan He, Yali Wang,
 Limin Wang, and Yu Qiao. VideoMamba: State space
 model for efficient video understanding. *arXiv preprint arXiv:2403.06977*, 2024. 1, 2
- [2] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco
 Massa, Alexandre Sablayrolles, and Hervé Jégou. Training
 data-efficient image transformers & distillation through attention. In *ICML*, pages 10347–10357. PMLR, 2021. 2