DOES VECTOR QUANTIZATION FAIL IN SPATIO-TEMPORAL FORECASTING? EXPLORING A DIFFER-ENTIABLE SPARSE SOFT-VECTOR QUANTIZATION AP-PROACH

Anonymous authors

008

009

010 011 012

013

015

016

017

018

019

021

022

024

025

026

027

028

029

031

Paper under double-blind review

ABSTRACT

Spatio-temporal forecasting is crucial in various fields and requires a careful balance between identifying subtle patterns and filtering out noise. Vector quantization (VQ) appears well-suited for this purpose, as it quantizes input vectors into a set of codebook vectors or patterns. Although VQ has shown promise in various computer vision tasks, it surprisingly falls short in enhancing the accuracy of spatio-temporal forecasting. We attribute this to two main issues: inaccurate optimization due to non-differentiability and limited representation power in hard VQ. To tackle these challenges, we introduce Differentiable Sparse Soft-Vector Quantization (SVQ), the first VQ method to enhance spatio-temporal forecasting. SVQ balances detail preservation with noise reduction, offering full differentiability and a solid foundation in sparse regression. Our approach employs a two-layer MLP and an extensive codebook to streamline the sparse regression process, significantly cutting computational costs while simplifying training and improving performance. Empirical studies on five spatio-temporal benchmark datasets show SVQ achieves state-of-the-art results, including a 7.9% improvement on the WeatherBench-S temperature dataset and an average MAE reduction of 9.4% in video prediction benchmarks (Human3.6M, KTH, and KittiCaltech), along with a 17.3% enhancement in image quality (LPIPS). Code is publicly available at https://anonymous.4open.science/r/SVQ-Forecasting.

- 033 034 1 IN
 - 1 INTRODUCTION

Spatio-temporal forecasting is pivotal in numerous domains ranging from environmental monitoring to urban planning, where precisely predicting future dynamics is crucial. The journey to refine 037 forecasting methods has spanned from traditional feature engineering to the latest explorations in deep learning. Among various methodologies explored, Vector Quantization (VQ) has distinguished itself primarily in computer vision tasks, showcasing its ability to compress high-dimensional 040 vectors into a compact, discrete form that maintains significant fidelity to the original information. 041 Historically rooted in signal processing, VQ's breakthrough came with its application in image 042 processing advancements like the Vector Quantised-Variational AutoEncoder (VQ-VAE) van den Oord et al. (2017), which set a precedent in generating high-quality images by learning efficient 043 representations of complex distributions. 044

045 While VQ has proven effective and has become a nearly default approach in computer vision 046 generation tasks, its potential applications in spatio-temporal forecasting remain less explored. 047 Considering the noise reduction capabilities of VQ, along with the similarities between image/video 048 generation and spatio-temporal forecasting, one could infer that VQ would positively impact the latter. 049 However, our review of existing studies reveals that few have successfully enhanced spatio-temporal forecasting performance using VQ techniques. Our empirical analysis of recent state-of-the-art VQ 050 methods discovered that all of them fell short of expectations, often degrading the performance 051 of baseline forecasting models rather than providing the anticipated improvements, as illustrated 052 in Figure 1 and detailed in Table 4. They consistently exert a significant negative influence on the final MSE or MAE regression accuracy.

054We believe that the similarity between image/video055generation and spatio-temporal forecasting under-056scores the potential of VQ, but the problem lies in the057inherent dynamic nature of spatio-temporal data. The058complex temporal evolutions and spatial distributions059introduce a level of complexity that traditional VQ060methods are not equipped to handle. Specifically, the061unsatisfactory results of traditional VQ methods are062caused by two reasons:

Inaccurate model optimization caused by non differentiability. The discrete nature of the quan tization step prevents gradients from being directly
 passed through this operation. VQ methods typically



Figure 1: Limitations of VQ in spatio-temporal forecasting: An experiment study evaluating mean squared error (MSE) improvement percentage on the WeatherBench-S temperature dataset.

employ the straight-through (or stop-gradient) estimator, as described in VQVAE van den Oord et al.
 (2017). This estimator approximates the gradient by copying gradients from the quantized outputs to the input vectors, introducing errors in the optimization process.

Limited representation power of hard-VQ. VQ methods typically assign each input vector to a
 single nearest codebook vector, which limits the modeling of detailed spatio-temporal dynamics
 required for forecasting.

In response to these limitations, this work presents Differentiable Sparse Soft-Vector Quantization (SVQ), a novel technique designed to strike a balance between noise reduction and detail preservation for spatio-temporal forecasting tasks. We solve the aforementioned challenges by:

- Introducing a differentiable VQ that simplifies gradient computations. This is achieved by approximating sparse regression with an MLP layer, coupled with a codebook, retaining the differentiability crucial for modern deep learning pipelines. A two-layer MLP generates regression coefficients through nonlinear projections of input vectors. The quantized outputs are derived from the dot product of these coefficients and the codebook matrix. Since the coefficients are generated from input vectors, gradients can flow directly from the quantized outputs to the input vectors. This straightforward yet effective approach not only enables differentiable VQ, enhancing accuracy, but also addressing the computational challenges often associated with sparse regression in VQ.
 - Using Soft-VQ with sparse regression to combine vectors from a large codebook. As shown
 in Figure 2, SVQ innovatively integrates sparse regression and allows for the allocation of
 input vectors with multiple codebook vectors. This significantly enhances the model's ability
 to capture intricate patterns and effectively filter out noise. Compared to hard-VQ, SVQ
 exhibits a more uniform distribution of codebook vectors, indicating that SVQ is able to
 preserve more diverse and fine-grained information from the original inputs. Our empirical
 studies reveal that SVQ possesses intriguing properties, such as effectively utilizing a
 completely frozen, randomly initialized codebook without sacrificing performance, thereby
 significantly reducing learning parameters and showcasing its efficiency and robustness.

0

076

077

078

079

081

082

084

085

090

092





107 Through rigorous testing on a variety of real-world datasets, **SVQ has proven to be the first VQ** method to achieve significant enhancements in spatio-temporal forecasting tasks. Notably, SVQ surpassed the leading model in the WeatherBench-S temperature forecasting benchmark by 7.9%. In
video prediction tasks—Human3.6M, KTH, and KittiCaltech, SVQ systematically lowered the Mean
Absolute Error (MAE) by 9.4%, while also marking a significant improvement in perceptual quality,
as indicated by a 17.3% reduction in the LPIPS score. These results underscore SVQ's remarkable
capability across a wide range of spatio-temporal forecasting tasks.

114 2 RELATED WORK

113

Due to space limitations, here we provide a brief overview of vector quantization, the lineage of sparse coding techniques, and the latest developments in spatio-temporal forecasting algorithms. An extensive review can be found in the Appendix B.

119 Vector Quantization and Sparse Coding. Instead of using continuous latent, VQ-VAE van den 120 Oord et al. (2017), a seminal work, incorporates vector quantization to learn discrete latent represen-121 tations, typically assigning each vector to the nearest code in a codebook. Subsequent enhancements 122 include Residual VQ Zeghidour et al. (2022), which quantizes the residuals recursively, and Multi-123 headed VO Mama et al. (2021b), which adopts multiple heads for each vector. While these methods 124 are effective, they often rely on a relatively small number of codes to represent the original vectors. 125 To address this, SCVAE Xiao et al. (2023) employs sparse coding, allowing vectors to be represented through sparse linear combinations of multiple codes, and achieves end-to-end training via the 126 Learnable Iterative Shrinkage Thresholding Algorithm (LISTA) Gregor & LeCun (2010). However, a 127 significant drawback of the sparse coding method using LISTA (referred to as SVQ-raw here) is its 128 high computational complexity, which scales quadratically with codebook size. 129

Building on these insights and limitations, our work proposes a soft-VQ method applicable to spatiotemporal forecasting tasks. Although it is closely related to a simultaneous research work Tschannen et al. (2023), which employs an infinite cookbook with a linear layer for continuous vector quantization applied in image generation, our approach is largely inspired by sparse regression, as clearly evidenced by our analysis. Specifically, our work focuses on the challenges arising from spatio-temporal forecasting, providing a strong theoretical foundation and effectively addressing the challenges.

136 Spatio-Temporal Forecasting Models. Recent advancements in spatio-temporal forecasting have 137 highlighted a shift from recurrent to non-recurrent frameworks. Despite the forecasting capabilities of 138 models like ConvLSTM SHI et al. (2015), PredRNN Wang et al. (2017), and PredRNNV2Wang et al. (2022), this shift is largely due to the high computational demands of sequential processing in recurrent 139 models. Non-recurrent models, such as MMVP Zhong et al. (2023) and the SimVP family Gao et al. 140 (2022); Tan et al. (2022), have become benchmarks in video prediction by decoupling spatial and 141 temporal learning through an efficient encoder-translator-decoder structure. This transition is further 142 enhanced by innovative features like visual attention in TAU Tan et al. (2023a) and MetaFormers 143 in OpenSTL Tan et al. (2023b), showcasing the continuous improvements towards more effective 144 forecasting solutions. Our proposed method is designed for seamless integration as a plugin with the 145 majority of these spatio-temporal forecasting models. 146

147 148

149

150

151

152

153 154

155

3 DIFFERENTIABLE SPARSE SOFT-VECTOR QUANTIZATION (SVQ)

In this section, we will first outline the mathematical foundation of sparse soft-vector quantization, followed by a detailed implementation within a spatio-temporal forecasting model. Our proposed method effectively addresses the **optimization problem** through differentiation, and the subsequent theoretical analysis of cookbook utilization demonstrates its substantial **representational capacity**.

3.1 VECTOR QUANTIZATION BY SPARSE REGRESSION

Let $\{z_i \in \mathbb{R}^d\}_{i=1}^m$ be the set of codes. A typical vector quantization method assigns a data point $x \in \mathbb{R}^d$ to the nearest code in $\{z_i\}_{i=1}^m$. The main problem with such an approach is that a significant part of the information in x will be lost due to quantization. Sparse regression turns the code assignment problem into an optimization problem

160 161 $w = \underset{w \in \mathbb{R}^{m}_{+}}{\operatorname{arg\,min}} \frac{1}{2} \left| x - \sum_{i=1}^{m} w_{i} z_{i} \right|^{2} + \lambda |w|_{1}, \tag{1}$ where $w = (w_1, \ldots, w_m) \in \mathbb{R}^m_+$ is the weight for combining codes $\{z_i\}_{i=1}^m$ to approximate x. λ refers to the regularization parameter. By introducing L_1 regularizer in the optimization problem, we effectively enforce x to be associated with a small number of codes. Compared to classic VQ methods where codes have to be learned through clustering, according to Chiu et al. (2022), it is sufficient to use randomly sampled vectors as codes as long as its number is large enough, thus avoiding the need of computing and adjusting codes. The theoretical guarantee of sparse regression is closely related to the property of subspace clustering, as revealed in Theorem 4.1.

As shown in Figure 3, the obvious downside of sparse regression for VQ is its high computational cost, as it needs to solve the optimization problem in (1) for EVERY data point. Below, we will show that sparse regression can be approximated by a two-layer MLP and a randomly fixed or learnable matrix, making it computationally attractive.

To solve the optimization problem (1), we consider the composite optimization method whose iteration is given as follows

$$w_{t+1}' = w_t - \eta Z^{\top} (Zw_t - x), \qquad (2)$$

$$[w_{t+1}] = \operatorname{sgn} \left([w_{t-1}'] \right) \left([[w_{t-1}']] - \lambda n \right) \qquad (3)$$

$$[w_{t+1}]_i = \operatorname{sgn}\left(\left[w_{t+1}'\right]_i\right)\left(\left|\left[w_{t+1}'\right]_i\right| - \lambda\eta\right)_+, (3)$$

182 where $Z = (z_1, \ldots, z_m)$, $[z]_i$ is the *i*th element 183 of vector *z*, sgn denotes the sign function, and $(a)_+$ 184 outputs 0 if a < 0 and *a* otherwise. We con-185 sider the first step of the iteration where $w_0 = 0$ 186 and have $w = \eta \text{sgn} (Z^T x - \lambda \mathbf{1}) [Z^T x - \lambda \mathbf{1}]_+ =$ 187 $\eta \text{sgn} (Z^T x - \lambda) \sigma (Z^T x - \lambda)$ and the resulting output 188 for *x* is given as 189



Figure 3: Effect of SVQ approximation: Floating point operations per second (FLOPS) and mean squared error (MSE) on WeatherBench-S temperature dataset with SVQ-raw and SVQ. The computational complexity of SVQ-raw increases quadratically with the size of codebook, making it suffer from out-of-memory (OOM) issue when scaling codebook size up to 2^{12} .

178

179

181

191 192 $x' = \sum_{i=1}^{m} w_i z_i = \eta Z sgn\left(Z^T x - \lambda\right) \sigma\left(Z^T x - \lambda\right).$ (4)

By generalizing ηZ into another matrix B, we have output vector x' exactly expressed as a matrix and a two-layer MLP over x. We finally note that although it is convenient to form the codebook by randomly sampling vectors, we found empirically that tuning codebook does bring slight additional gains in some cases.

197

199

3.2 SPATIO-TEMPORAL FORECASTING MODEL ENHANCED BY QUANTIZATION

200 Architecture of backbone model. As depicted in Figure 4, SimVP Tan et al. (2022) is employed as 201 the backbone model, which encompasses an encoder for spatial feature extraction, a translator for temporal dependency learning, and a decoder for frame reconstruction. The quantization module is 202 integrated between the encoder and translator. The input data is a 4D tensor $X \in \mathbb{R}^{H*W*T*C}$, repre-203 senting height (H), width (W), time step (T), and channel (C). The encoder En condenses X into 204 downsampled latent representation $En(X) \in \mathbb{R}^{H'*W'*T*C'}$, maintaining temporal dimensionality 205 while altering spatial and channel dimensions. This latent space, composed of H' * W' * T tokens, 206 each represented by a C'-dimensional vector, undergoes vector quantization. 207

208 Quantization module. The SVQ comprises a two-layer MLP and an extensive codebook. The 209 codebook is a randomly initialized matrix $\mathcal{M} \in \mathbb{R}^{N * C'}$, where N denotes the size of codebook. 210 To achieve automatic selection of codes, a weight matrix $\mathcal{W} \in \mathbb{R}^{H'*W'*T*N}$ is generated via 211 nonlinear projection from the latent representation En(X). This projection is formally expressed 212 as $\mathcal{W} = MLP(En(X))$, wherein the MLP comprises two linear layers and an intermediate ReLU 213 activation function. The quantized output Q is then obtained by computing the dot product of weight matrix \mathcal{W} and codebook matrix \mathcal{M} , a process that can be conceptualized as a selection operation 214 as shown in Figure 4. To encourage sparsity within the generated weight matrix, we apply a Mean 215 Absolute Error (MAE) loss to the output as a surrogate form of regularization.



Figure 4: **Top:** Architecture of backbone model and the proposed quantization module. The encoder, translator, decoder are inherited from SimVP. A quantization module is added between the encoder and translator to effectively ensure a good generalized performance. Bottom: Quantization process of traditional VQ (Left) and our proposed SVQ (Right). In contrast, SVQ select multiple codes (red dots) from a huge codebook (gray dots), and the codebook can be either learnable or frozen.

4 EFFICIENT UTILIZE OF COOKBOOK USING SPARSE REGRESSION

To understand the difference between sparse regression-based quantization scheme and clusteringbased quantization scheme, we measure the number of codes required to approximate any vector within a unit ball \mathcal{B} with error less than δ . This number is denoted by $T(\mathcal{B}, \delta)$. Intriguingly, as the theorem below reveals, using sparse regression allows $T(\mathcal{B}, \delta)$ to be significantly reduced from $O(1/\delta^d)$ to $O(1/\delta^p)$, where $p \ll d$ for high-dimensional vectors.

Theorem 4.1. For the clustering-based method, $T(\mathcal{B}, \delta)$ is at least $1/\delta^d$. In contrast, for sparse regression, $T(\mathcal{B}, \delta)$ can be formulated as $(4d/\delta)^p$, where

$$p \ge \max\left(3, \frac{\log(4/\delta)}{\log\log(2d/\varepsilon)}\right),$$
(5)

given that the number of non-zero elements utilized by sparse regression is at least

$$\frac{4d}{\delta\left(\log C + p\log(4d) - (p+1)\log\delta\right)}.$$
(6)

Proof. To estimate $T(\mathcal{B}, \delta)$ for the clustering method, we consider the covering number for a unit 252 ball \mathcal{B} which necessitates at least $1/\delta^d$ code vectors to approximate any vector within an acceptable 253 error margin of δ . With $U = (u_1, \ldots, u_m)$ where $u_k \sim \mathcal{N}(0, I_d/m)$, and $g \in \Delta_s$ an s-sparse unit 254 vector, we discern:

$$\Pr\left(\|UU^{\top} - I\|_{2} \ge \gamma\right) \le 2d \exp\left(-\frac{m\gamma^{2}}{3d}\right),\tag{7}$$

which implies

$$\|UU^{\top} - I\|_{2} \le \Delta := \sqrt{\frac{d}{m} \log \frac{2d}{\varepsilon}}$$
(8)

with probability at least $1 - \varepsilon$. Therefore, $\|g' - g\|_2 \ge (1 + \Delta)^{-1} \|Ug - Ug'\|_2$. Since the *s*-sparse unit vector covering number is bounded by $(Cm/s\delta)^s$, we establish:

$$\left(\frac{Cm}{s\delta}\right)^s \ge \left(1 + \frac{2}{\delta}\right)^d (1 + \Delta)^d,\tag{9}$$

Setting $m = (4d/\delta)^p$ yields

$$(1+\Delta)^d \le \exp(d\Delta) \le e,\tag{10}$$

therefore, $s \log(C'm/\delta) \ge d(2/\delta + \log(1 + \Delta))$, where C' = Ce. As long as $s \ge \frac{4d}{\delta(\log C + p \log(4d) - (p+1)\log \delta)}$, it follows that $s \log s \le 2d/\delta$, affirming that $m \ge (4d/\delta)^p$.

²⁷⁰ 5 EXPERIMENTS

271 272

We extensively evaluate SVQ on a wide range of real-world spatio-temporal forecasting datasets under the unified framework of OpenSTL Tan et al. (2023b). Given that SimVP holds leading performance across almost all benchmarks, it serves as our primary baseline.

Dataset. We conduct extensive experiments on five real-world spatio-temporal forecasting tasks, 276 including weather (WeatherBench Rasp et al. (2020)), traffic flow (TaxiBJ Zhang et al. (2017)), human 277 pose dynamics (Human3.6M Ionescu et al. (2014)), driving scenes (KittiCaltech Geiger et al. (2013); 278 Dollár et al. (2009)), and human actions (KTH Action Schüldt et al. (2004)). The above datasets have 279 relatively few channels. As the number of channels increases, it becomes more challenging to apply 280 VQ, requiring a more diverse codebook. Therefore, to validate the performance on high-dimensional 281 data, additional experiments were conducted using the WeatherBench dataset in a High-dimensional 282 Multi-Variable (HMV) setting, which includes a total of 110 meteorological factors. Details about 283 datasets are provided in Appendix A.1.

Experimental details. During deployment, we found SVQ to be quite robust to codebook size, as its performance remains consistently strong when using a sufficiently large codebook. Therefore, we fix the codebook size at 10,000 for WeatherBench, TaxiBJ, and Human3.6M datasets, and at 6,000 for KittiCaltech and KTH datasets. The hidden dimension of nonlinear projection layer is fixed at 128. Experiments are conducted on either 1 or 4 NVIDIA V100 32GB GPUs, with a total batch size of 16. More details about backbone architectures, VQ parameters, computational costs, and metrics are described in Appendix A.2, A.3, D.1, and A.5, respectively.

291 292

293

5.1 BENCHMARKS ON VARIOUS FORECASTING TASKS

We explore both fixed (frozen) and learnable versions of SVQ on various forecasting tasks. Interest-295 ingly, our findings reveal that with a large codebook size, the performance of a frozen, randomlyinitialized codebook is on par with that of a carefully learned codebook. This observation aligns 296 with our intuition: when allowed to choose a very large number of representative vectors to form 297 a codebook, a random choice is often as good as the one that is carefully chosen, which has al-298 ready been studied in the column subset selection problem in matrix theory Drineas et al. (2008); 299 Deshpande & Rademacher (2010). The comparison baselines consist of two categories: 1) Non-300 recurrent models including SimVP Tan et al. (2022) and TAU Tan et al. (2023a); 2) Recurrent-based 301 models including ConvLSTM SHI et al. (2015), PredNet Lotter et al. (2017), PredRNN Wang et al. 302 (2017), PredRNN++ Wang et al. (2018), MIM Wang et al. (2019b), E3D-LSTM Wang et al. (2019a), 303 PhyDNet Guen & Thome (2020), MAU Chang et al. (2021), PredRNNv2 Wang et al. (2022), and 304 DMVFN Hu et al. (2023). Baseline results are copied from the original OpenSTL paper Tan et al. 305 (2023b). To preclude ambiguity, we select the best MetaFormer of SimVP for each dataset, detailed 306 in Appendix A.2.

The benchmark results of WeatherBench and three video prediction datasets (Human3.6M, KTH, and KittiCaltech) are presented in Tables 1 and 2, respectively. Due to page limit, results of WeatherBench-HMV and TaxiBJ datasets are provided in Appendix D.8 and D.7. These datasets have different characteristics. WeatherBench and TaxiBJ are macro forecasting tasks with low-frequency collection (30min or 1-6h). Human3.6M features subtle, low-frequency frame differences. KittiCaltech is challenging due to rapidly changing backgrounds and limited training data. The KTH dataset tests long-horizon forecasting, requiring the prediction of 20 future frames from 10 observed frames.

314 However, despite the distinct characteristics among datasets, a common thread is the need for 315 improved noise reduction coupled with enhanced representational capabilities, which can universally 316 benefit their respective forecasting tasks. Notably, the SimVP+SVQ model achieves either the best 317 or comparable performance across all datasets. For instance, on the WeatherBench-S temperature 318 dataset, SVQ significantly improves the best baseline by 7.9% (1.105 \rightarrow 1.018). On these three popular video prediction tasks, SVQ not only delivers a reduction in forecasting errors (average 9.4% 319 decrease in MAE), but also significantly improves subjective image quality (average 17.3% decrease 320 in LPIPS). On the WeatherBench-HMV dataset, SVQ continues to demonstrate a reduction in MAEs 321 in 110 channels, with an average of 8.9%. The results affirm that SVQ maintains good performance 322 when applied to high-dimensional datasets. Additional visualizations of forecasting samples can be 323 found in Appendix F.

Table 1: WeatherBench results: Performance comparison for SVQ module and baseline models on Weather-Bench. WeatherBench-S is single-variable, one-hour interval forecasting setup trained on data from 2010-2015, validated on 2016, and tested on 2017-2018. WeatherBench-M targets broader application, which is multivariable, six-hour interval forecasting setup trained on data from 1979-2015, validated on 2016, and tested on 2017-2018. The best and the second best results are highlighted by **bold** and underlined.

Dataset	Variable	Tempo	erature	Hum	idity	Wind C	omponent	Total Cl	oud Cover
	Model		MAE↓	MSE↓	MAE↓	MSE↓	MAE↓	MSE↓	MAE↓
	ConvLSTMSHI et al. (2015)	1.521	0.7949	35.146	4.012	1.8976	0.9215	0.0494	0.1542
	E3D-LSTMWang et al. (2019a)	1.592	0.8059	36.534	4.100	2.4111	1.0342	0.0573	0.1529
	PredRNNWang et al. (2017)	1.331	0.7246	37.611	4.096	1.8810	0.9068	0.0550	0.1588
	MIMWang et al. (2019b)	1.784	0.8716	36.534	4.100	3.1399	1.1837	0.0573	0.1529
WeatherBench-S	MAUChang et al. (2021)	1.251	0.7036	34.529	4.004	1.9001	0.9194	0.0496	0.1516
Weather Denen-5	PredRNN++Wang et al. (2018)	1.634	0.7883	35.146	4.012	1.8727	0.9019	0.0547	0.1543
	PredRNN.V2Wang et al. (2022)	1.545	0.7986	36.508	4.087	2.0072	0.9413	0.0505	0.1587
	TAUTan et al. (2023a)	1.162	0.6707	31.831	3.818	1.5925	0.8426	0.0472	0.1460
	SimVP (w/o VQ)Tan et al. (2022)	1.105	0.6567	31.332	3.776	1.4996	0.8145	0.0466	0.1469
	SimVP+SVQ (Frozen codebook)	<u>1.023</u>	<u>0.6131</u>	<u>30.863</u>	3.661	<u>1.4337</u>	<u>0.7861</u>	0.0456	0.1456
	SimVP+SVQ (Learnable codebook)	1.018	0.6109	30.611	3.657	1.4186	0.7858	<u>0.0458</u>	0.1463
	Improvement	↑7.9%	↑7.0%	↑2.3%	↑3.2%	↑5.4%	↑3.5%	↑ 2.1%	↑0.9%
	Variable	Tempe	erature	Hum	idity	Wind U G	Component	Wind V	Component
	ConvLSTMSHI et al. (2015)	6.303	1.7695	368.15	13.490	30.002	3.8923	30.789	3.8238
	PredRNNWang et al. (2017)	5.596	1.6411	354.57	13.169	27.484	3.6776	28.973	3.6617
	MIMWang et al. (2019b)	7.515	1.9650	408.24	14.658	35.586	4.2842	36.464	4.2066
WeatherBench-M	MAUChang et al. (2021)	5.628	1.6810	363.36	13.503	27.582	3.7409	27.929	3.6700
	PredRNN++Wang et al. (2018)	5.647	1.6433	363.15	13.246	28.396	3.7322	29.872	3.7067
	PredRNN.V2Wang et al. (2022)	6.307	1.7770	368.52	13.594	29.833	3.8870	31.119	3.8406
	TAUTan et al. (2023a)	4.904	1.5341	<u>342.63</u>	12.801	24.719	3.5060	25.456	3.4723
	SimVP (w/o VQ)Tan et al. (2022)	4.833	1.5246	340.06	12.738	24.535	3.4882	25.232	3.4509
	SimVP+SVQ (Frozen codebook)	4.427	1.4160	360.15	12.445	<u>23.915</u>	3.4078	24.968	<u>3.4117</u>
	SimVP+SVQ (Learnable codebook)	4.433	<u>1.4164</u>	360.53	<u>12.449</u>	23.908	3.4060	<u>24.983</u>	3.4095
	Improvement	↑8.4%	↑7.1%	↓5.9%	↑2.3%	↑2.6%	↑2.4%	↑1.0%	↑1.2%
	Dataset WeatherBench-S WeatherBench-M	Dataset Variable Model Model ConvLSTMSHI et al. (2015) E3D-LSTMWang et al. (2019a) PredRNNWang et al. (2017) MIMWang et al. (2017) MMWang et al. (2019b) MAUChang et al. (2017) PredRNNVang et al. (2017) MIMWang et al. (2018) PredRNN++Wang et al. (2018) PredRNN++Wang et al. (2022) TAUTan et al. (2023a) SimVP +SVQ (Frozen codebook) SimVP+SVQ (Learnable codebook) Improvement Variable ConvLSTMSHI et al. (2015) PredRNNNWang et al. (2017) MIMWang et al. (2017) MIMWang et al. (2019b) MAUChang et al. (2017) PredRNN++Wang et al. (2019) PredRNN++Wang et al. (2018) PredRNNN-V2Wang et al. (2018) PredRNN++Wang et al. (2018) PredRNN-V2Wang et al. (2022) TAUTan et al. (2023a) SimVP (w/o VQ)Tan et al. (2022) SimVP+SVQ (Frozen codebook) SimVP+SVQ (Learnable codebook) SimVP+SVQ (Frozen codebook)	Dataset Variable Temp Model MSEJ GonvLSTMSHI et al. (2015) 1.521 E3D-LSTMWang et al. (2019a) 1.592 PredRNNWang et al. (2017) 1.331 MIMWang et al. (2019b) 1.784 MAUChang et al. (2017) 1.331 MIMWang et al. (2017b) 1.784 MAUChang et al. (2017b) 1.634 PredRNN+4Wang et al. (2018) 1.634 PredRNN.V2Wang et al. (2022) 1.162 SimVP (w/o VQ)Tan et al. (2022) 1.105 SimVP+SVQ (Frozen codebook) 1.023 SimVP+SVQ (Learnable codebook) 1.023 Variable Temp ConvLSTMSHI et al. (2017) 5.596 MAUChang et al. (2017) 5.628 PredRNN-Y2Wang et al. (2017) 5.628 PredRNN-Y2Wang et al. (2017) 5.628 PredRNN-Y2Wang et al. (2021) 5.628 PredRNN-Y2Wang et al. (2022) 6.307 TredRNN-Y2Wang et al. (2022) 6.307 PredRNN-Y2Wang et al. (2022) 4.394 SimVP+SVQ (Frozen codebook) 4.433	Dataset Variable Temperature Model MSE↓ MAE↓ Model MSE↓ MAE↓ ConvLSTMSHI et al. (2015) 1.521 0.7949 E3D-LSTMWang et al. (2019a) 1.592 0.8059 PredRNNWang et al. (2017) 1.331 0.7246 MIMWang et al. (2019b) 1.784 0.8716 MAUChang et al. (2021) 1.251 0.7036 PredRNN+twang et al. (2022) 1.545 0.7986 TAUTan et al. (2023a) 1.162 0.6707 SimVP+sVQ (Frozen codebook) 1.023 0.6131 SimVP+sVQ (Learnable codebook) 1.023 0.6131 Io18 0.6109 1.018 0.6109 Yariable Temperature ConvLSTMSHI et al. (2015) 5.596 1.6411 MIMWang et al. (2017) 5.598 1.6421 1.6423 1.6411 MIMWang et al. (2017) 5.628 1.6810 1.6433 1.46431 PredRNN-Wang et al. (2018) 5.647 1.6433 1.4643 1.5246 SimVP+SVQ (Frozen codebook) </th <th>Dataset Variable Temperature Hum Model MSE↓ MAE↓ MSE↓ MSE↓ MSE↓ WeatherBench-S ConvLSTMSHI et al. (2019a) 1.521 0.7949 35.146 PredRNNWang et al. (2019a) 1.592 0.8059 36.534 PredRNNWang et al. (2017) 1.331 0.7246 37.611 MIMWang et al. (2019b) 1.784 0.8716 36.534 MAUChang et al. (2021) 1.251 0.7036 34.529 PredRNN++Wang et al. (2018) 1.644 0.7883 35.146 PredRNN-V2Wang et al. (2022) 1.545 0.7986 36.508 TAUTan et al. (2023a) 1.162 0.6707 31.831 SimVP+SVQ (Frozen codebook) 1.0023 0.6131 30.863 SimVP+SVQ (Frozen codebook) 1.0018 0.6109 30.611 Mathwang et al. (2015) 6.303 1.7695 368.15 PredRNNWang et al. (2017) 5.596 1.6411 354.57 MIMWang et al. (2019) 7.515 1.9650 408.54.57</th> <th>Dataset Variable Temperature Humidity Model MSE↓ MAE↓ MSE↓ MAE↓ ConvLSTMSHI et al. (2015) 1.521 0.7949 35.146 4.012 E3D-LSTMWang et al. (2019a) 1.592 0.8059 36.534 4.100 PredRNNWang et al. (2017) 1.331 0.7246 37.611 4.096 MIMWang et al. (2019b) 1.784 0.8716 36.534 4.100 MAUChang et al. (2018) 1.634 0.7883 35.146 4.012 PredRNN+4Wang et al. (2012) 1.545 0.7986 36.508 4.007 TAUTan et al. (2023a) 1.162 0.6707 31.831 3.818 SimVP+sVQ (Frozen codebook) 1.003 0.6617 30.863 3.661 SimVP+SVQ (Frozen codebook) 1.018 0.6109 30.863 3.661 SimVP+SVQ (Learnable codebook) 1.018 0.6109 30.6513 3.490 WeatherBench-M Variable Temperature Humidity ConvLSTMSHI et al. (2015) 6.303 1.7695</th> <th>Dataset Variable Temperature Humidity Wind C Model MSE↓ MAE↓ MSE↓ MAE↓ MSE↓ MAE↓ MSE↓ WeatherBench-S ConvLSTMSHI et al. (2015) 1.521 0.7949 35.146 4.012 1.8976 WeatherBench-S E3D-LSTMWang et al. (2017) 1.331 0.7246 37.611 4.096 1.8810 MIMWang et al. (2019b) 1.784 0.8716 36.534 4.100 3.1399 MAUChang et al. (2019b) 1.251 0.7036 34.529 4.004 1.9001 PredRNN++Wang et al. (2018) 1.654 0.7883 35.146 4.0121 1.8727 PredRNN+5VQ (Frozen codebook) 1.022 1.545 0.7986 36.508 4.002 1.8727 SimVP+SVQ (Frozen codebook) 1.023 0.6131 3.818 1.5925 1.4996 SimVP+SVQ (Frozen codebook) 1.023 0.6131 3.657 1.4186 Improvement ↑7.9% ↑7.0% ↑2.3% ↑3.2% ↑5.4% Weath</th> <th>Dataset Variable Temperature Humidity Wind Component Model MSE↓ MAE↓ MSE↓ MAE↓ MSE↓ MAE↓ MAE↓</th> <th>Dataset Variable Temperature Humidity Wind Component Total Cl Model MSE↓ MAE↓ MSE↓ MAE↓ MSE↓ MAE↓ MAE↓ MSE↓ MAE↓ MAE↓ MSE↓ 0.0494 Bab_LSTMWang et al. (2019) 1.592 0.8059 36.534 4.100 2.4111 1.0342 0.0573 MAUChang et al. (2021) 1.251 0.7036 34.529 4.004 1.9001 0.9194 0.0496 PredRNN++Wang et al. (2022) 1.254 0.7986 36.508 4.072 1.8727 0.9019 0.0547 TAUTan et al. (2023a) 1.162 0.6707 31.831 3.818 1.5925 0.8426 0.0472 SimVP+SVQ (Frozen codebook) 1.018 0.06131 36.651 1.4186 0.7858 0.0458 </th>	Dataset Variable Temperature Hum Model MSE↓ MAE↓ MSE↓ MSE↓ MSE↓ WeatherBench-S ConvLSTMSHI et al. (2019a) 1.521 0.7949 35.146 PredRNNWang et al. (2019a) 1.592 0.8059 36.534 PredRNNWang et al. (2017) 1.331 0.7246 37.611 MIMWang et al. (2019b) 1.784 0.8716 36.534 MAUChang et al. (2021) 1.251 0.7036 34.529 PredRNN++Wang et al. (2018) 1.644 0.7883 35.146 PredRNN-V2Wang et al. (2022) 1.545 0.7986 36.508 TAUTan et al. (2023a) 1.162 0.6707 31.831 SimVP+SVQ (Frozen codebook) 1.0023 0.6131 30.863 SimVP+SVQ (Frozen codebook) 1.0018 0.6109 30.611 Mathwang et al. (2015) 6.303 1.7695 368.15 PredRNNWang et al. (2017) 5.596 1.6411 354.57 MIMWang et al. (2019) 7.515 1.9650 408.54.57	Dataset Variable Temperature Humidity Model MSE↓ MAE↓ MSE↓ MAE↓ ConvLSTMSHI et al. (2015) 1.521 0.7949 35.146 4.012 E3D-LSTMWang et al. (2019a) 1.592 0.8059 36.534 4.100 PredRNNWang et al. (2017) 1.331 0.7246 37.611 4.096 MIMWang et al. (2019b) 1.784 0.8716 36.534 4.100 MAUChang et al. (2018) 1.634 0.7883 35.146 4.012 PredRNN+4Wang et al. (2012) 1.545 0.7986 36.508 4.007 TAUTan et al. (2023a) 1.162 0.6707 31.831 3.818 SimVP+sVQ (Frozen codebook) 1.003 0.6617 30.863 3.661 SimVP+SVQ (Frozen codebook) 1.018 0.6109 30.863 3.661 SimVP+SVQ (Learnable codebook) 1.018 0.6109 30.6513 3.490 WeatherBench-M Variable Temperature Humidity ConvLSTMSHI et al. (2015) 6.303 1.7695	Dataset Variable Temperature Humidity Wind C Model MSE↓ MAE↓ MSE↓ MAE↓ MSE↓ MAE↓ MSE↓ WeatherBench-S ConvLSTMSHI et al. (2015) 1.521 0.7949 35.146 4.012 1.8976 WeatherBench-S E3D-LSTMWang et al. (2017) 1.331 0.7246 37.611 4.096 1.8810 MIMWang et al. (2019b) 1.784 0.8716 36.534 4.100 3.1399 MAUChang et al. (2019b) 1.251 0.7036 34.529 4.004 1.9001 PredRNN++Wang et al. (2018) 1.654 0.7883 35.146 4.0121 1.8727 PredRNN+5VQ (Frozen codebook) 1.022 1.545 0.7986 36.508 4.002 1.8727 SimVP+SVQ (Frozen codebook) 1.023 0.6131 3.818 1.5925 1.4996 SimVP+SVQ (Frozen codebook) 1.023 0.6131 3.657 1.4186 Improvement ↑7.9% ↑7.0% ↑2.3% ↑3.2% ↑5.4% Weath	Dataset Variable Temperature Humidity Wind Component Model MSE↓ MAE↓ MSE↓ MAE↓ MSE↓ MAE↓ MAE↓	Dataset Variable Temperature Humidity Wind Component Total Cl Model MSE↓ MAE↓ MSE↓ MAE↓ MSE↓ MAE↓ MAE↓ MSE↓ MAE↓ MAE↓ MSE↓ 0.0494 Bab_LSTMWang et al. (2019) 1.592 0.8059 36.534 4.100 2.4111 1.0342 0.0573 MAUChang et al. (2021) 1.251 0.7036 34.529 4.004 1.9001 0.9194 0.0496 PredRNN++Wang et al. (2022) 1.254 0.7986 36.508 4.072 1.8727 0.9019 0.0547 TAUTan et al. (2023a) 1.162 0.6707 31.831 3.818 1.5925 0.8426 0.0472 SimVP+SVQ (Frozen codebook) 1.018 0.06131 36.651 1.4186 0.7858 0.0458

> > Table 2: Video prediction results: Performance comparison for SVQ module and baseline models on Human3.6M, KTH, and KittiCaltech. The best and the second best results are highlighted by **bold** and underlined.

Dataset	Dataset Human3.6M			KittiCaltech				KTH				
Metric	MAE↓	SSIM↑	PSNR ↑	LPIPS↓	MAE↓	SSIM↑	PSNR↑	LPIPS↓	MAE↓	SSIM↑	PSNR↑	LPIPS↓
ConvLSTMSHI et al. (2015)	1583.3	0.9813	33.40	0.03557	1583.3	0.9345	27.46	0.08575	445.5	0.8977	26.99	0.26686
E3D-LSTMWang et al. (2019a)	1442.5	0.9803	32.52	0.04133	1946.2	0.9047	25.45	0.12602	892.7	0.8153	21.78	0.48358
PredNetLotter et al. (2017)	1625.3	0.9786	31.76	0.03264	1568.9	0.9286	27.21	0.11289	783.1	0.8094	22.45	0.32159
PhyDNetGuen & Thome (2020)	1614.7	0.9804	39.84	0.03709	2754.8	0.8615	23.26	0.32194	765.6	0.8322	23.41	0.50155
MAUChang et al. (2021)	1577.0	0.9812	33.33	0.03561	1800.4	0.9176	26.14	0.09673	471.2	0.8945	26.73	0.25442
MIMWang et al. (2019b)	1467.1	0.9829	33.97	0.03338	1464.0	0.9409	28.10	0.06353	380.8	0.9025	27.78	0.18808
PredRNNWang et al. (2017)	1458.3	0.9831	33.94	0.03245	1525.5	0.9374	27.81	0.07395	380.6	0.9097	27.95	0.21892
PredRNN++Wang et al. (2018)	1452.2	0.9832	34.02	0.03196	1453.2	0.9433	28.02	0.13210	370.4	0.9124	28.13	0.19871
PredRNN.V2Wang et al. (2022)	1484.7	0.9827	33.84	0.03334	1610.5	0.9330	27.12	0.08920	368.8	0.9099	28.01	0.21478
TAUTan et al. (2023a)	1390.7	0.9839	34.03	0.02783	1507.8	0.9456	27.83	0.05494	421.7	0.9086	27.10	0.22856
DMVFN Hu et al. (2023)	-	-	-	-	1531.1	0.9314	26.95	0.04942	413.2	0.8976	26.65	0.12842
SimVP (w/o VQ)Tan et al. (2022)	1441.0	0.9834	34.08	0.03224	1507.7	0.9453	27.89	0.05740	397.1	0.9065	27.46	0.26496
SimVP+SVQ (Frozen codebook)	1264.9	0.9851	34.07	0.02380	1408.6	0.9469	28.10	0.05535	364.6	0.9109	27.28	0.20988
SimVP+SVQ (Learnable codebook)	1265.1	0.9851	34.06	0.02367	1414.9	0.9458	28.10	0.05776	360.2	0.9116	27.37	0.20658
Improvement	↑12 .2%	↑0.2%	$\downarrow 0.0\%$	↑26.2%	↑6.6%	↑0. 2%	↑0.8%	↑3.6%	↑9.3%	↑ 0.6%	↓0.3%	↑22.0%

5.2 BOOSTING PERFORMANCE AS A VERSATILE PLUG-IN

In this section, SVQ serves as a versa- Table 3: Boosting performance: The effect of tile plug-in module applicable to various MetaFormers Yu et al. (2022a). The adopted MetaFormers include three types. 1) CNN-based: SimVPv1(IncepU) Gao et al. (2022), SimVPv2(gSTA) Tan et al. (2022), Con-vMixer Trockman & Kolter (2023), Con-vNeXt Liu et al. (2022b), HorNet Rao et al. (2022), and MogaNet Li et al. (2022). 2) Transformer-based: ViT Dosovitskiy et al. (2021), Swin Transformer Liu et al. (2021), Uni-former Li et al. (2023), Poolformer Yu et al. (2022b), and VAN Guo et al. (2023). 3) MLP-based: MLPMixer Tolstikhin et al. (2021). We conduct experiments on WeatherBench-S temperature dataset because it is lightweight and fast for training.

SVQ for various MetaFormers on WeatherBench-S temperature dataset.

MetaFormer	MS	Е	MAE		
	w/o SVQ	w SVQ	w/o SVQ	w SVQ	
SimVPv1(IncepU)Gao et al. (2022)	1.238	1.216	0.7037	0.6831	
SimVPv2(gSTA)Tan et al. (2022)	1.105	1.018	0.6567	0.6109	
ConvMixerTrockman & Kolter (2023)	1.267	1.257	0.7073	0.6780	
ConvNeXtLiu et al. (2022b)	1.277	1.159	0.7220	0.6568	
HorNetRao et al. (2022)	1.201	1.130	0.6906	0.6472	
MogaNetLi et al. (2022)	1.152	1.067	0.6665	0.6271	
ViTDosovitskiy et al. (2021)	1.146	1.111	0.6712	0.6375	
SwinLiu et al. (2021)	1.143	1.088	0.6735	0.6320	
UniformerLi et al. (2023)	1.204	1.110	0.6885	0.6400	
PoolformerYu et al. (2022b)	1.156	1.097	0.6715	0.6297	
VANGuo et al. (2023)	1.150	1.083	0.6803	0.6342	
MLP-MixerTolstikhin et al. (2021)	1.255	1.120	0.7011	0.6455	
Average improvement	↑ 4.8	8%	↑ 6.0%		

378 As shown in Table 3, SVQ consistently improves the performance across all MetaFormers, showcasing 379 its universality across diverse backbone types. We observe an average reduction in MSE and MAE 380 by 4.8% and 6.0%, respectively. In detail, SVQ leads to an average MSE reduction of 4.1%381 for CNN-based backbones, 5.1% for transformer-based, and 10.7% for MLP-based. The more 382 pronounced enhancement in transformer-based and MLP-based models indicates that our approach is especially effective with architectures that prioritize global interactions. Notably, SimVPv2(gSTA) 383 is the best backbone, while our SVQ further improves it by 7.9%. These findings also aligns with 384 our motivation that mitigating noise in the learning process significantly benefits spatio-temporal 385 forecasting, irrespective of model architecture. By integrating SVQ to constrain the diversity of 386 predicted patterns and cut out noise, researchers can focus on crafting high-quality and general base 387 models. 388

389 390

5.3 DELICATE BALANCE BETWEEN DETAIL PRESERVATION AND NOISE REDUCTION

391 To study the role of VQ in spatio-temporal fore-392 casting, we evaluated several cutting-edge VQ 393 methods akin to the SVQ framework, imple-394 mented as plug-in modules alongside the back-395 bone forecasting model. Table 4 shows that 396 SVQ significantly improves forecasting as a 397 plug-in, whereas other VQ methods result in 398 increased prediction errors. Enhanced detail retention within the representational capacity is 399 associated with lower forecasting errors. Clas-400 sic VO methods suffer from notable information 401 losses, as evidenced by a higher MSE of 1.854. 402 In contrast, both residual VQ and grouped resid-403 ual VQ outperform traditional VQ with lower 404

Table 4: Comparison of vector quantization methods: All methods share identical backbone, with the recommended setting in Appendix A.3. The results better than baseline are highlighted in **bold**.

Method	MSE↓	MAE↓
Baseline (SimVP w/o VQ)	1.105	0.6567
VQ van den Oord et al. (2017)	1.854	0.8963
Residual VQ (RVQ) Zeghidour et al. (2022)	1.213	0.6910
Grouped Residual VQ (GRVQ) Yang et al. (2023)	1.174	0.6747
Multi-headed VQ (MHVQ) Mama et al. (2021b)	1.211	0.6994
Stochastic Residual VQ (SRVQ) Lee et al. (2022)	1.888	0.9237
Residual Finite Scalar Quantization (RFSQ) Mama et al. (2021a)	1.319	0.7505
Lookup Free Quantization (LFQ) Yu et al. (2023a)	2.988	1.1103
Residual LFQ (RLFQ) Yu et al. (2023a)	1.281	0.7281
SVQ-raw Xiao et al. (2023)	1.123	0.6456
SVQ	1.018	0.6109

MSEs of 1.213 and 1.174, respectively, affirming their ability to preserve intricate details due to 405 recursive quantization.

406

407 408 It is commonly understood that the codebook size in 409 clustering-based VQ is critical: larger codebooks cap-MSE 410 ture more details, whereas smaller ones enhance noise 1.4 411 reduction. To explore this trade-off, we compared how the Prediction 1.2 412 codebook size influences the prediction MSE in Grouped 413 Residual VQ (GRVQ) and SVQ. As Figure 5 indicates, the MSE of GRVQ initially decreases but increases with 414 overly large codebooks, echoing findings from Yu et al. 415 (2023b) that an excessively large codebook may degrade 416 image generation performance. This underscores the ne-417 cessity for dataset-specific tuning in clustering-based VQ 418

approaches. In contrast, SVQ naturally achieves a balance



Figure 5: Predition MSE curves on WeatherBench-S temperature dataset with Grouped Residual VQ (GRVQ) and SVQ.

419 between preserving detail and reducing noise through sparse regression, eliminating the need for 420 extensive fine-tuning. The codebook size in SVQ exhibits a low-maintenance profile: using a default 421 large codebook can produce robust results without extensive tuning. We are not suggesting that our 422 SVQ outperforms other VQ methods in general image generation tasks, as it is beyond the scope 423 of our current objective. Rather, we emphasize SVQ's effectiveness as a noise reduction tool that 424 directly enhances real-world spatio-temporal forecasting tasks, while the application to general image generation remains a topic for future exploration. 425

426

427 5.4 TRAIN STABILITY ISSUE 428

429 Although it is feasible to place the quantization module either before or after the translator, we found that for post-translator placement, the traditional VQ method van den Oord et al. (2017) suffers 430 pronounced instability and codebook collapse issues, as shown in Figure 6. It is essential to highlight 431 that the backbones without VQ maintain their MSE within the acceptable range of approximately 1 to

2 (refer to Table 3). Yet, integrating traditional VQ causes a substantial rise in MSE values, exceeding 10 for different backbones—a level considered excessively high. We hypothesize that this instability is attributed to the non-differentiability of the straight-through estimator, which introduces errors into the gradient flow for preceding modules. In contrast, our SVQ module never encounters such issues and remains highly stable throughout training. To maintain the integrity of all VQ methods, we opt for the pre-translator design in our main experiments, wherein quantization is executed preceding the translator module. The difference between two designs is detailed in Appendix D.3.



Figure 6: VQ (Top) and SVQ (Bottom) training curves: We perform post-translator quantization on various backbones, with the same codebook size (1024), employing early stopping (patience of 10) on the WeatherBench-S temperature dataset. Perplexity for SVQ is averaged over different θ values, detailed in Appendix A.4.

5.5 ABLATION STUDY

We conduct a series of ablation studies on WeatherBench-S temperature dataset to understand the contribution of important designs based on the default setting: SVQ with a codebook size of 10,000, learnable codebook, and MAE loss. An additional ablation study on the frozen module is provided in Appendix D.5.

Self-learned sparse regression structure. The original SimVP model adopts MSE as prediction loss. We individually replace it with MAE loss and add the SVQ module. As shown in Table 5, the joint use of SVQ and MAE loss is crucial for significantly improving the model's performance. We suggest that the sparsity of the weight matrix \mathcal{W} impacts vector representation learning and use kurtosis to quantify this after normalizing \mathcal{W} . Figure 7 demonstrates that both a learnable codebook and MAE loss contribute to increased sparsity. Analyzing four codebook initialization methods in both learnable and fixed settings (Table 6), we find that a learnable codebook promotes sparsity irrespective of the initialization, indicating that sparsity is a self-learned property that enhances intermediate representation learning.

Table 5:	Ablation	of SVQ	Compoments
----------	----------	--------	------------

481	Module	MSE↓	MAE↓
482	SimVP (MSE loss)	1.105	0.6567
483	SimVP (MAE loss)	1.126	0.6509
100	SimVP+SVQ (Learnable, MSE loss)	1.099	0.6527
484	SimVP+SVQ (Learnable, MAE loss)	1.018	0.6109
485			

Initialization	Learnability	MSE↓	MAE↓	Kurtosis
kaiming uniform	Frozen	1.023	0.6131	1.596
	Learnable	1.018	0.6109	7.213
sparse(sparsity=0.9)	Frozen	1.050	0.6183	4.165
	Learnable	1.034	0.6160	41.558
trunc normal	Frozen	1.049	0.6166	1.582
	Learnable	1.031	0.6161	4.236
orthogonal	Frozen	1.034	0.6170	1.561
	Learnable	1.030	0.6131	35.774



Figure 7: Distribution of regression weight \mathcal{W} and codebook \mathcal{M} : Higher kurtosis represents more compact and concentrate distribution near zero, as well as sparser regression weights. Left: Learnable SVQ with MAE loss. Middle: Learnable SVQ with MSE loss. Right: Frozen SVQ with MAE loss. Learnable setting and MAE loss encourage sparser weights and a more structured codebook.

502 Codebook size and learnability. Table 7 503 compares the effects of codebook size-both 504 learnable and frozen. Results show that 505 increasing codebook size consistently en-506 hances performance. However, when the 507 size reaches 10,000, the performance gap be-508 tween frozen and learnable codebooks nar-509 rows to just 0.5% (1.023 \rightarrow 1.018). A larger 510 codebook provides comprehensive coverage of the latent space through random codes, 511 minimizing the need for meticulous learning. 512 Consequently, models with randomly initial-513

Table 7: Ablation of Model Structure.

Projection dim	Codebook size	MSE↓	MAE↓
128	10	1.070	0.6227
128	1000	1.044	0.6198
128	10000	1.023	0.6131
128	10	1.060	0.6194
128	1000	1.048	0.6182
128	10000	1.018	0.6109
1280 (Bucket-shape)	1280	1.035	0.6149
None (One-layer)	10000	1.043	0.6144
128 (Post-ReLU)	10000	1.032	0.6136
	Projection dim 128 128 128 128 128 128 128 128	Projection dim Codebook size 128 10 128 1000 128 10000 128 1000 128 1000 128 1000 128 1000 128 1000 128 10000 1280 (Bucket-shape) 1280 None (One-layer) 10000 128 (Post-ReLU) 10000	Projection dim Codebook size MSE↓ 128 10 1.070 128 1000 1.044 128 10000 1.023 128 10000 1.044 128 1000 1.048 128 1000 1.048 128 1000 1.048 128 10000 1.048 1280 (Bucket-shape) 1280 1.035 None (One-layer) 10000 1.043 128 (Post-ReLU) 10000 1.032

ized codebooks perform similarly to those with learned ones. Additionally, our optimized SVQ 514 structure outperforms alternative designs, including single-layer, bucket-shaped, and post-ReLU 515 variants, while keeping a similar parameter count. 516

5.6 ROBUSTNESS TO NOISE, ERROR-BAR, CONVERGENCE BEHAVIOUR, VISUALIZATIONS OF PREDICTIONS AND LATENT VECTORS, AND HIGH-DIMENSIONAL BENCHMARK RESULTS.

We conducted additional experiments by introducing artificial noise to the training data, confirming that our method effectively mitigates noise by constraining latent patterns through quantization, as detailed in Appendix D.2. The statistical significance of the error bars is provided in Appendix D.6. We further analyzed the convergence behavior of SVQ and traditional VQ in Appendix C. Additionally, we included supplementary experiments to understand the impact of SVQ on latent representation and to compare various VQ methods in Appendices E, and D.4, respectively. A benchmark experiment 526 for high-dimensional spatio-temporal forecasting is also included in Appendix D.8.

527 528 529

517 518

519

520 521

522

523

524

498

499

500

501

CONCLUSIONS 6

530 531

532 In this work, we present Differentiable Sparse Soft-Vector Quantization (SVQ), a concise yet effective 533 method for spatio-temporal forecasting enhancement. Unlike other state-of-the-art VQ methods, this 534 is the first approach that demonstrates a boosting effect in spatio-temporal forecasting tasks.SVQ elegantly tackles the inaccuracies in the optimization problem arising from non-differentiability and the restricted representational capabilities associated with hard-VQ. Tested across diverse benchmarks, 537 from weather to traffic and video prediction, SVQ consistently outperforms pure baseline methods, setting new performance standards without complex priors. Its differentiability and seamless inte-538 gration with baseline models highlight SVQ as a significant advancement for efficient and effective spatio-temporal forecasting.

540 REFERENCES 541

567

568

569

581

- Zheng Chang, Xinfeng Zhang, Shanshe Wang, Siwei Ma, Yan Ye, Xiang Xinguang, and Wen 542 Gao. MAU: A motion-aware unit for video prediction and beyond. In Marc'Aurelio Ran-543 zato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan 544 (eds.), Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 546 26950-26962, 2021. URL https://proceedings.neurips.cc/paper/2021/hash/ 547 e25cfa90f04351958216f97e3efdabe9-Abstract.html. 548
- Chung-Cheng Chiu, James Qin, Yu Zhang, Jiahui Yu, and Yonghui Wu. Self-supervised learning 549 with random-projection quantizer for speech recognition. In International Conference on Machine 550 Learning, 2022. 551
- 552 Amit Deshpande and Luis Rademacher. Efficient volume sampling for row/column subset selection. 553 In 51th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2010, October 554 23-26, 2010, Las Vegas, Nevada, USA, pp. 329-338. IEEE Computer Society, 2010. doi: 10.1109/ FOCS.2010.38. URL https://doi.org/10.1109/FOCS.2010.38. 555
- 556 Piotr Dollár, Christian Wojek, Bernt Schiele, and Pietro Perona. Pedestrian detection: A benchmark. In 2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 558 2009), 20-25 June 2009, Miami, Florida, USA, pp. 304-311. IEEE Computer Society, 2009. doi: 10. 559 1109/CVPR.2009.5206631. URL https://doi.org/10.1109/CVPR.2009.5206631.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 561 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, 562 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. 563 In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. URL https://openreview.net/forum?id= 565 YicbFdNTTy. 566
 - Petros Drineas, Michael W. Mahoney, and S. Muthukrishnan. Relative-error CUR matrix decompositions. SIAM J. Matrix Anal. Appl., 30(2):844-881, 2008. doi: 10.1137/07070471X. URL https://doi.org/10.1137/07070471X.
- 570 Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image 571 synthesis. In IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, 572 virtual, June 19-25, 2021, pp. 12873-12883. Computer Vision Foundation / IEEE, 2021. doi: 10.1109/CVPR46437.2021.01268. URL https://openaccess.thecvf.com/content/ 573 CVPR2021/html/Esser_Taming_Transformers_for_High-Resolution_ 574 Image_Synthesis_CVPR_2021_paper.html. 575
- 576 Zhangyang Gao, Cheng Tan, Lirong Wu, and Stan Z. Li. Simvp: Simpler yet better video prediction. 577 In IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, 578 LA, USA, June 18-24, 2022, pp. 3160–3170. IEEE, 2022. doi: 10.1109/CVPR52688.2022.00317. 579 URL https://doi.org/10.1109/CVPR52688.2022.00317.
- 580 Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The KITTI dataset. Int. J. Robotics Res., 32(11):1231–1237, 2013. doi: 10.1177/0278364913491297. 582 URL https://doi.org/10.1177/0278364913491297.
- 584 Karol Gregor and Yann LeCun. Learning fast approximations of sparse coding. In Johannes Fürnkranz and Thorsten Joachims (eds.), Proceedings of the 27th International Conference on 585 Machine Learning (ICML-10), June 21-24, 2010, Haifa, Israel, pp. 399–406. Omnipress, 2010. 586 URL https://icml.cc/Conferences/2010/papers/449.pdf.
- 588 Vincent Le Guen and Nicolas Thome. Disentangling physical dynamics from unknown factors 589 for unsupervised video prediction. In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, pp. 11471-11481. Computer Vision Foundation / IEEE, 2020. doi: 10.1109/CVPR42600.2020.01149. URL https://openaccess.thecvf.com/content_CVPR_2020/html/Le_ 592 Guen_Disentangling_Physical_Dynamics_From_Unknown_Factors_for_ Unsupervised_Video_Prediction_CVPR_2020_paper.html.

617

627

632

633

634

635

636

637

638

639

640

641

- 594
 595
 595
 596
 596
 597
 Meng-Hao Guo, Cheng-Ze Lu, Zheng-Ning Liu, Ming-Ming Cheng, and Shi-Min Hu. Visual attention network. *Comput. Vis. Media*, 9(4):733–752, 2023. doi: 10.1007/S41095-023-0364-2.
 597
- Xiaotao Hu, Zhewei Huang, Ailin Huang, Jun Xu, and Shuchang Zhou. A dynamic multi-scale voxel flow network for video prediction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pp. 6121–6131. IEEE, 2023. doi: 10.1109/CVPR52729.2023.00593. URL https://doi.org/10.1109/CVPR52729.2023.00593.
- Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Trans. Pattern Anal. Mach. Intell.*, 36(7):1325–1339, 2014. doi: 10.1109/TPAMI.2013.248. URL https://doi.org/10.1109/TPAMI.2013.248.
- Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Autoregressive image
 generation using residual quantization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp. 11513–11522. IEEE, 2022.
 doi: 10.1109/CVPR52688.2022.01123. URL https://doi.org/10.1109/CVPR52688.
 2022.01123.
- Kunchang Li, Yali Wang, Junhao Zhang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. Uniformer: Unifying convolution and self-attention for visual recognition. *IEEE Trans. Pattern Anal. Mach. Intell.*, 45(10):12581–12600, 2023. doi: 10.1109/TPAMI.2023.3282631. URL https://doi.org/10.1109/TPAMI.2023.3282631.
- Siyuan Li, Zedong Wang, Zicheng Liu, Cheng Tan, Haitao Lin, Di Wu, Zhiyuan Chen, Jiangbin Zheng, and Stan Z. Li. Efficient multi-order gated aggregation network. *CoRR*, abs/2211.03295, 2022. doi: 10.48550/ARXIV.2211.03295. URL https://doi.org/10.48550/arXiv.2211.03295.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In 2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021, pp. 9992–10002. IEEE, 2021. doi: 10.1109/ICCV48922.2021.00986. URL https: //doi.org/10.1109/ICCV48922.2021.00986.
- Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp. 3192–3201. IEEE, 2022a. doi: 10.1109/CVPR52688. 2022.00320. URL https://doi.org/10.1109/CVPR52688.2022.00320.
 - Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, *CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pp. 11966–11976. IEEE, 2022b. doi: 10. 1109/CVPR52688.2022.01167. URL https://doi.org/10.1109/CVPR52688.2022. 01167.
 - William Lotter, Gabriel Kreiman, and David D. Cox. Deep predictive coding networks for video prediction and unsupervised learning. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net, 2017. URL https://openreview.net/forum?id=B1ewdt9xe.
- Rayhane Mama, Marc S. Tyndel, Hashiam Kadhim, Cole Clifford, and Ragavan Thurairatnam. NWT: towards natural audio-to-video generation with representation learning. *CoRR*, abs/2106.04283, 2021a. URL https://arxiv.org/abs/2106.04283.
- Rayhane Mama, Marc S. Tyndel, Hashiam Kadhim, Cole Clifford, and Ragavan Thurairatnam. NWT: towards natural audio-to-video generation with representation learning. *CoRR*, abs/2106.04283, 2021b. URL https://arxiv.org/abs/2106.04283.

680

681

682

683

684

685

686

687

 Yongming Rao, Wenliang Zhao, Yansong Tang, Jie Zhou, Ser-Nam Lim, and Jiwen Lu. Hornet: Efficient high-order spatial interactions with recursive gated convolutions. In *NeurIPS*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/ 436d042b2dd81214d23ae43eb196b146-Abstract-Conference.html.

- Stephan Rasp, Peter D Dueben, Sebastian Scher, Jonathan A Weyn, Soukayna Mouatadid, and Nils
 Thuerey. Weatherbench: a benchmark data set for data-driven weather forecasting. *Journal of Advances in Modeling Earth Systems*, 12(11):e2020MS002203, 2020.
- Christian Schüldt, Ivan Laptev, and Barbara Caputo. Recognizing human actions: A local SVM approach. In *17th International Conference on Pattern Recognition, ICPR 2004, Cambridge, UK, August 23-26, 2004*, pp. 32–36. IEEE Computer Society, 2004. doi: 10.1109/ICPR.2004.1334462.
 URL https://doi.org/10.1109/ICPR.2004.1334462.
- Xingjian SHI, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-kin Wong, and Wangchun WOO. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper/2015/file/ 07563a3fe3bbe7e3ba84431ad9d055af-Paper.pdf.
- Cheng Tan, Zhangyang Gao, and Stan Z. Li. Simvp: Towards simple yet powerful spatiotemporal
 predictive learning. *CoRR*, abs/2211.12509, 2022. doi: 10.48550/ARXIV.2211.12509. URL
 https://doi.org/10.48550/arXiv.2211.12509.
- 670 Cheng Tan, Zhangyang Gao, Lirong Wu, Yongjie Xu, Jun Xia, Siyuan Li, and Stan Z. Li. Temporal attention unit: Towards efficient spatiotemporal predictive learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pp. 18770–18782. IEEE, 2023a. doi: 10.1109/CVPR52729.2023.01800. URL https://doi.org/10.1109/CVPR52729.2023.01800.
- 675
 676
 676
 676
 677
 677
 678
 678
 Cheng Tan, Siyuan Li, Zhangyang Gao, Wenfei Guan, Zedong Wang, Zicheng Liu, Lirong Wu, and Stan Z. Li. Openstl: A comprehensive benchmark of spatio-temporal predictive learning. *CoRR*, abs/2306.11249, 2023b. doi: 10.48550/ARXIV.2306.11249. URL https://doi.org/10.48550/arXiv.2306.11249.
 - Ilya O. Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and Alexey Dosovitskiy. Mlp-mixer: An all-mlp architecture for vision. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 24261–24272, 2021. URL https://proceedings.neurips.cc/paper/2021/hash/cba0a4ee5ccd02fda0fe3f9a3e7b89fe-Abstract.html.
- Asher Trockman and J. Zico Kolter. Patches are all you need? *Trans. Mach. Learn. Res.*, 2023, 2023. URL https://openreview.net/forum?id=rAnB7JSMXL.
- Michael Tschannen, Cian Eastwood, and Fabian Mentzer. Givt: Generative infinite-vocabulary transformers. ArXiv, abs/2312.02116, 2023. URL https://api.semanticscholar.org/ CorpusID:265610025.
- Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp. 6306–6315, 2017. URL https://proceedings.neurips.cc/ paper/2017/hash/7a98af17e63a0ac09ce2e96d03992fbc-Abstract.html.
- Yunbo Wang, Mingsheng Long, Jianmin Wang, Zhifeng Gao, and Philip S. Yu. Predrnn: Recurrent neural networks for predictive learning using spatiotemporal lstms. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman

702	Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference
703	on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA,
705	pp. 8/9-888, 2017. UKL https://proceedings.neurips.cc/paper/2017/nash/
706	estoadoces/41//eetozsbisdocorobo-Abscract.nemi.
707	Yunbo Wang, Zhifeng Gao, Mingsheng Long, Jianmin Wang, and Philip S. Yu. Predrnn++: Towards
708	A resolution of the deep-in-time dilemma in spatiotemporal predictive learning. In Jennifer G.
709	Dy and Andreas Krause (eds.), Proceedings of the 35th International Conference on Machine
710	Learning, ICML 2018, Stockholmsmassan, Stockholm, Sweden, July 10-15, 2018, volume 80 of Ducceedings of Machine Learning Percent pp 5110 5110 DMLP 2018, UPL http://
711	//proceedings mlr press/w80/wang18b html
712	//proceedings.mir.press/voo/wangrob.nemr.
713	Yunbo Wang, Lu Jiang, Ming-Hsuan Yang, Li-Jia Li, Mingsheng Long, and Li Fei-Fei. Eidetic 3d
714	LSTM: A model for video prediction and beyond. In 7th International Conference on Learning
715	Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019a.
716	URL https://openreview.net/forum?id=BIIKS2AqtX.
717	Yunbo Wang, Jianjin Zhang, Hongyu Zhu, Mingsheng Long, Jianmin Wang, and Philip S. Yu.
718	Memory in memory: A predictive neural network for learning higher-order non-stationarity from
719	spatiotemporal dynamics. In IEEE Conference on Computer Vision and Pattern Recognition, CVPR
720	2019, Long Beach, CA, USA, June 16-20, 2019, pp. 9154–9162. Computer Vision Foundation
721	/ IEEE, 2019D. doi: 10.1109/CVPK.2019.0095/. UKL http://openaccess.thecvi.
722	Neural Network for Learning Higher-Order CVPR 2019 paper html
723	Neurar_Neework_ror_hearning_higher order_evin_zory_paper.neur.
724	Yunbo Wang, Haixu Wu, Jianjin Zhang, Zhifeng Gao, Jianmin Wang, Philip Yu, and Mingsheng Long.
725	Predrnn: A recurrent neural network for spatiotemporal predictive learning. <i>IEEE Transactions on</i>
720	Pattern Analysis and Machine Intelligence, pp. 1–1, 2022. doi: 10.1109/1PAMI.2022.3165153.
728	Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment:
729	from error visibility to structural similarity. IEEE Trans. Image Process., 13(4):600-612, 2004.
730	doi: 10.1109/TIP.2003.819861. URL https://doi.org/10.1109/TIP.2003.819861.
731	Pan Xiao Peijie Oju, and Aristeidis Sotiras, SC-VAE: sparse coding-based variational autoencoder
732	<i>CoRR</i> , abs/2303.16666, 2023. doi: 10.48550/ARXIV.2303.16666. URL https://doi.org/
733	10.48550/arXiv.2303.16666.
734	
735	Ziru Xu, Yunbo Wang, Mingsheng Long, and Jianmin Wang. Predcnn: Predictive learning with
736	Laint Conference on Artificial Intelligence IICAI 2018 July 13-19 2018 Stockholm Sweden
737	pp. 2940–2947. jicai.org. 2018. doi: 10.24963/IJCAL2018/408. URL https://doi.org/10.
738	24963/ijcai.2018/408.
739	
740	Dongenao rang, Songxiang Liu, Kongjie Huang, Jinenuan Iian, Chao Weng, and Yuexian Zou. Hif-
741	2023 doi: 10.48550/ARXIV2305.02765 URL https://doi.org/10.48550/arXiv
742	2305.02765.
743	
744	Lijun Yu, José Lezama, Nitesh Bharadwaj Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen,
746	Yong Uneng, Agrim Gupta, Xiuye Gu, Alexander G. Hauptmann, Boqing Gong, Ming-Hsuan Vang, Irfan Essa, David A. Poss, and Lu Jiang. Language model basts diffusion to basis in large
747	to visual generation CoRR abs/2310.05737, 2023a, doi: 10.48550/ARXIV.2310.05737, UPI
748	https://doi.org/10.48550/arXiv.2310.05737
749	
750	Lijun Yu, José Lezama, Nitesh Bharadwaj Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen,
751	Yong Uneng, Agrim Gupta, Xiuye Gu, Alexander G. Hauptmann, Boqing Gong, Ming-Hsuan
752	to visual generation CoRR abs/2310.05737 2023b. doi: 10.48550/ARXIV.2310.05737 UPI
753	https://doi.org/10.48550/arXiv.2310.05737
754	,, ac,, atht. doi:0.00.0
755	Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In <i>IEEE/CVF Conference</i>

on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pp. 10809-10819. IEEE, 2022a. doi: 10.1109/CVPR52688.2022.01055. URL https: //doi.org/10.1109/CVPR52688.2022.01055. Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pp. 10809-10819. IEEE, 2022b. doi: 10.1109/CVPR52688.2022.01055. URL https: //doi.org/10.1109/CVPR52688.2022.01055. Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. Soundstream: An end-to-end neural audio codec. IEEE ACM Trans. Audio Speech Lang. Process., 30: 495-507, 2022. doi: 10.1109/TASLP.2021.3129994. URL https://doi.org/10.1109/ TASLP.2021.3129994. Junbo Zhang, Yu Zheng, and Dekang Qi. Deep spatio-temporal residual networks for citywide crowd flows prediction. In Satinder Singh and Shaul Markovitch (eds.), Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA, pp. 1655–1661. AAAI Press, 2017. doi: 10.1609/AAAI.V3111.10735. URL https://doi.org/10.1609/aaai.v31i1.10735. Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018, pp. 586–595. Computer Vision Foundation / IEEE Computer Society, 2018. doi: 10.1109/CVPR. 2018.00068. URL http://openaccess.thecvf.com/content cvpr 2018/html/ Zhang_The_Unreasonable_Effectiveness_CVPR_2018_paper.html. Yiqi Zhong, Luming Liang, Ilya Zharkov, and Ulrich Neumann. MMVP: motion-matrix-based video prediction. In IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023, pp. 4250–4260. IEEE, 2023. doi: 10.1109/ICCV51070.2023.00394. URL https://doi.org/10.1109/ICCV51070.2023.00394.

Supplemental Materials

The supplementary material for our work *Does Vector Quantization Fail in Spatio-Temporal Forecasting? Exploring a Differentiable Sparse Soft-Vector Quantization Approach* is organized as follows: Appendix A provides implementation details of SimVP model and VQ methods. Appendix B gives an extensive review of related work. Appendix C analyzes the convergence behaviour of SVQ and traditional VQ. Appendix D present extended quantitative results. Appendix E delves deeper into the effect of SVQ on latent representation. Finally, Appendix F shows additional qualitative results of forecasting samples and errors.

819 820 821

822 823

824 825

827

828

829

842 843

844

845

846

847 848

810

811

A IMPLEMENTATION DETAILS

A.1 DATASET DETAILS

WeatherBench Rasp et al. (2020) and TaxiBJ Zhang et al. (2017) are two macro forecasting tasks collected at low frequencies (30min or 1-6h). Human3.6M Ionescu et al. (2014), KittiCaltech Geiger et al. (2013); Dollár et al. (2009), and KTH Action Schüldt et al. (2004) are three popular video prediction tasks. A summary of dataset statistics is provided in Table 8.

Table 8: The detailed statistics of benchmark datasets.

Dataset	Size train test		Seq. Len. in out		Img. Shape $H \times W \times C$	Interval
WeatherBench-S	2,167	706	12	12	$\begin{array}{c} 32 \times 64 \times 1 \\ 32 \times 64 \times 4 \\ 32 \times 32 \times 2 \\ 128 \times 160 \times 3 \\ 256 \times 256 \times 3 \end{array}$	1 hour
WeatherBench-M	54,019	2,883	4	4		6 hour
TaxiBJ	20,461	500	4	4		30 min
KittiCaltech	3,160	3,095	10	1		Frame
Human3.6M	73,404	8,582	4	4		Frame
KTH Action	4,940	3,030	10	20	$\begin{array}{c} 128 \times 128 \times 1 \\ 32 \times 64 \times 110 \end{array}$	Frame
WeatherBench-HMV	52,559	2,883	4	4		6 hour

A.2 ARCHITECTURE CONFIGURATION OF SIMVP

Table 9 reports the architectures of SimVP on all datasets. We select the best MetaFormer to replace the translator module based on OpenSTL benchmarks¹²³. The parameters remain unchanged, following the original configurations. It is noteworthy that due to reproducibility issues of ConvNeXt on the TaxiBJ dataset, we have opted to utilize gSTA as our backbone model.

Table 9: Detailed configuration of SimVP backbone.

Dataset MetaFormer (Translator)		spatio_kernel	hid_S	hid_T	N_T	N_S	drop_path	LR scheduler
WeatherBench-S temperature	gSTA	enc=3, dec=3	32	256	8	2	0.1	cosine
WeatherBench-S humidity	Swin	enc=3, dec=3	32	256	8	2	0.2	cosine
WeatherBench-S wind component	Swin	enc=3, dec=3	32	256	8	2	0.2	cosine
WeatherBench-S total cloud cover	gSTA	enc=3, dec=3	32	256	8	2	0.1	cosine
WeatherBench-M	MogaNet	enc=3, dec=3	32	256	8	2	0.1	cosine
TaxiBJ	gŠTA	enc=3, dec=3	32	256	8	2	0.1	cosine
Human3.6M	gSTA	enc=3, dec=3	64	512	6	4	0.1	cosine
KTH	IncepU	enc=3, dec=3	64	256	6	2	0.1	onecycle
KittiCaltech	gSTA	enc=3, dec=3	64	256	6	2	0.2	onecycle
WeatherBench-HMV	gSTA	enc=3, dec=3	32	256	8	2	0.1	cosine

857 858 859

861

862

863

¹https://openstl.readthedocs.io/en/latest/model_zoos/video_benchmarks. html

²https://openstl.readthedocs.io/en/latest/model_zoos/weather_benchmarks. html

³https://openstl.readthedocs.io/en/latest/model_zoos/traffic_benchmarks. html

A.3 PARAMETERS OF COMPARED VQ METHODS

Table 4 presents a comparison of SVQ with several well-known VQ methods, reproduced using source
code from the GitHub repository⁴. The parameters were kept consistent with the recommended
settings to ensure performance, as detailed in Table 10. It should be noted that we found that
increasing the codebook size for previous VQ methods, such as Residual VQ and Multi-headed VQ,
led to a considerable increase in GPU memory usage and extended the training time to impractical
levels. This issue is one of the reasons these methods recommend adopting a default codebook size
of 1024. To ensure fairness, we conducted an extensive experiment for VQ methods using the same
codebook size (1024) in Appendix D.4.

0	_	Л
0	ſ	4
8	7	5

Table 10: Parameters of the comp	pared VQ methods.
----------------------------------	-------------------

Vector quantization method	codebook_size	num_quantizers	groups	heads	shared_codebook	Specific parameters
VQ	512	-	-	-	-	-
Residual VQ	1024	8	-	-	\checkmark	-
Grouped Residual VQ	1024	8	2	-	\checkmark	-
Multi-headed VQ	1024	-	-	8	\checkmark	-
Residual VQ (Stochastic)	1024	8	-	-	\checkmark	stochastic_sample_codes=True
Residual Finite Scalar Quantization	-	8	-	-	-	levels=[8, 5, 5, 3]
Lookup Free Quantization (LFQ)	8192	-	-	-	-	entropy_loss_weight=0.1
Residual LFQ	256	8	-	-	-	

A.4 EVALUATION OF PERPLEXITY

Unlike other VQ methods that rely on a single code, our SVQ generates multiple regression weights to merge several codes. To evaluate its perplexity, we first normalize the regression weights and then convert them into binary form using a threshold set at θ times the standard deviation, where θ serves as the threshold value. We utilize two thresholds (2 and 3) to obtain reasonable perplexity.

A.5 METRICS

Forecasting accuracy is evaluated using mean squared error (MSE) and mean absolute error (MAE), while the image quality of predicted frames is assessed using structural similarity index measure (SSIM) Wang et al. (2004), peak signal-to-noise ratio (PSNR), and learned perceptual image patch similarity (LPIPS) Zhang et al. (2018). The training process is early stopped with a patience of 10, and the models with the minimal loss are saved for subsequent evaluation.

B EXTENSIVE REVIEW OF RELATED WORK

B.1 RECURRENT-BASED FORECASTING MODEL

The majority of spatio-temporal forecasting models leverage techniques such as Conv2D Xu et al. (2018), Conv3D Wang et al. (2019a), and attention mechanisms Liu et al. (2022a) for spatial modeling. Distinctions among these models primarily arise from how they incorporate temporal information. Recurrent-based models, exemplified by ConvLSTM SHI et al. (2015), have been widely used to capture motion dynamics by iteratively processing multi-frame predictions. Variants like PredRNN Wang et al. (2017) introduce the Spatio-Temporal LSTM (ST-LSTM) unit, integrating spatial appearances and temporal variations within a single memory pool. Further advancements include PredRNN++ Wang et al. (2018) and PredRNNV2 Wang et al. (2022), which deepen the model and expand the receptive field through a cascading LSTM mechanism and a memory decoupling strategy, respectively. MIM Wang et al. (2019b) network utilizes a self-renewed memory module to exploit differential signals by decomposing non-stationary dynamics. PredNet Lotter et al. (2017) improves performance by estimating prediction errors in forward propagation.

B.2 NON-RECURRENT FORECASTING MODEL

Despite the effectiveness of recurrent-based methods, they suffer from high computational cost caused
 by their inherent unparallelizable architecture. Recent efforts in spatio-temporal forecasting have
 shifted towards non-recurrent models that decouple the forecasting task from autoregressive processes.

⁴https://github.com/lucidrains/vector-quantize-pytorch/tree/master

918 Notably, SimVPv1 Gao et al. (2022) and SimVPv2 Tan et al. (2022) separate spatial and temporal 919 learning into distinct phases within an encoder-translator-decoder structure, consistently surpassing 920 recurrent counterparts in video prediction tasks. TAU Tan et al. (2023a) refines the architecture by 921 incorporating a visual attention mechanism into the translator. OpenSTL Tan et al. (2023b) further 922 enhances the translator by MetaFormers.

924 **B.3** VECTOR QUANTIZATION

926 By partitioning a continuous vector space into a discrete collection of vectors, VQ effectively reduces the data required to characterize a set of values, thereby achieving noise reduction. Following the 927 introduction of VO-VAE van den Oord et al. (2017), many variants such as VOGAN Esser et al. 928 (2021), Residual VQ Zeghidour et al. (2022), Multi-headed VQ Mama et al. (2021b), Grouped 929 Residual VQ Yang et al. (2023), Finite Scalar Quantization (FSQ) Mama et al. (2021a), and Lookup 930 Free Quantization Yu et al. (2023a) have been developed to enhance the representational capabilities 931 of VQ. For instance, FSQ Mama et al. (2021a) simplifies VQ in generative modeling by discretizing 932 scalar values. However, to the best of our knowledge, VQ has seen limited application in spatio-933 temporal forecasting, which has inspired our research. Traditional hard-VQ tends to eliminate 934 excessive detail, thus impairing forecasting accuracy. While the sparse coding-based variational 935 autoencoder (SC-VAE) Xiao et al. (2023) incorporates sparse coding into the variational autoencoder 936 framework, its application is primarily targeted at image reconstruction and segmentation tasks. Furthermore, in the main text, our experiments have demonstrated that the implementation of the 937 LISTA algorithm Gregor & LeCun (2010) within SC-VAE leads to out-of-memory issues when a 938 large codebook size is employed. 939

940 941

942

950

951

953 954 955

956

963 964 965

966 967

923

925

CONVERGENCE BEHAVIOR OF SVQ AND TRADITIONAL VQ С

943 SVQ and traditional VQ are based on sparse regression and clustering, respectively. Since SVQ is 944 fully differentiable, their convergence behaviors can be considered analogous to Backpropagation (BP) and K-Nearest Neighbors (KNN), respectively. We prove that BP's optimization is smoother 945 than KNN's due to continuous gradient descent. 946

947 **Backpropagation.** The decision function is directly related to the weights and activation functions, 948 which tend to evolve smoothly under gradient descent. 949

$$f(x) = \sigma(W_n(\sigma(W_{n-1}(\dots\sigma(W_1x + b_1)\dots) + b_{n-1})) + b_n),$$
(11)

with the gradient updates affecting W_i and b_i : 952

$$W_i \leftarrow W_i - \eta \nabla_{W_i} J,\tag{12}$$

$$b_i \leftarrow b_i - \eta \nabla_{b_i} J. \tag{13}$$

The gradient $\nabla L(\mathbf{W})$ is typically smooth and continuous if the loss function L and the activation 957 functions are smooth. Hence, the weights update in a relatively smooth manner, and the path to the 958 minimum of the loss function is traversed in small, continuous steps. 959

960 K-Nearest Neighbors. Being a non-parametric method, KNN directly relies on the training data to 961 make its predictions. For a new input \mathbf{x} , the prediction y is made based on the majority vote among its k nearest neighbors: 962

$$y = \operatorname{argmax} \sum_{i=1}^{k} \delta(y_i, y), \tag{14}$$

968 where δ is the Kronecker delta function. 969

The decision boundary of KNN is piecewise linear and can change abruptly with small changes in 970 the input data. Consider the case where the input \mathbf{x} moves slightly from one side of the decision 971 boundary to the other:

974
$$\mathbf{x} \to \mathbf{x}' | \mathbf{x} \approx \mathbf{x}'.$$
 (15)

In this case, the set of
$$k$$
 nearest neighbors might change abruptly, leading to a discontinuous jump in the prediction y .

We can conclude that the optimization process for backpropagation is smoother as compared to KNN due to the continuous nature of gradient descent, which updates the weights in small, incremental steps:

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla L(\mathbf{W}_t). \tag{16}$$

In contrast, KNN's decision boundary can lead to discontinuous predictions with small perturbations in the input data, illustrating the non-smooth nature of the optimization process in KNN.

D ADDITIONAL QUANTITATIVE RESULTS

D 1 COMPUTATIONAL COST

Table 11 presents the computational costs of SVQ module and forecasting models. It shows that recurrent-based models have significantly higher FLOPS requirements, while non-recurrent models are more efficient. The proposed SVQ module is not only effective but also computationally cheap. Across all datasets, SVQ only slightly adds the number of parameters and FLOPS. The computational burden of SimVP+SVQ remains significantly smaller than recurrent-based models.

Table 11: Number of parameters and computing performances for all forecasting models.

1004	Model type	Dataset	Huma	m3.6M	К	TH	KittiC	altech	Weather	Bench-S	Tax	iBJ
1005		Model	Params	FLOPS	Params	FLOPS	Params	FLOPS	Params	FLOPS	Params	FLOPS
1005		ConvLSTM	15.5M	347.0G	14.9M	1368.0G	15.0M	595.0G	14.98M	136G	14.98M	20.74G
1007		PredNet	12.5M	542.0G 13.7G	12.5M	217.0G 3.4G	12.5M	1004G 12.5M	51.09M	169G -	12.5M	98.19G 0.85G
1007	D (1 1	PhyDNet	4.2M	19.1G	3.1M	93.6G	3.1M	40.4G	3.09M	36.8G	3.09M	5.60G
1008	Recurrent-based	MAU MIM	20.2M 47.6M	105.0G 1051.0G	20.1M 39.8M	399.0G 1099.0G	24.3M 49.2M	172.0G 1858G	37.75M	39.6G 109G	4.41M 37.86M	6.02G 64.10G
1009		PredRNN	24.6M	704.0G	23.6M	2800.0G	23.7M	1216G	23.57M	278G	23.66M	42.40G
1010		PredRNN.V2	24.6M	708.0G	23.6M	2815.0G	23.8M	1223G	23.59M	279G	23.67M	42.63G
1011		DMVFN	-	-	3.5M	0.88G	3.6M	1.2G	-	-	3.54M	0.057G
1012	Non-recurrent	TAU SimVP (w/o VQ)	37.6M 28.8M	182.0G 146.0G	15.0M 12.2M	73.8G 62.8G	44.7M 15.6M	80.0G 96.3G	12.22M 12.76M	6.70G 7.01G	9.55M 7.84M	2.49G 2.08G
1013		SimVP+SVQ	30.7M	178.0G	13.3M	110.0G	16.8M	156G	14.37M	16.8G	9.45M	3.72G

D.2 ROBUSTNESS TO ARTIFICIAL NOISE INJECTION

To clearly demonstrate the noise reduction effect of SVQ, we conduct a series of experiments by introducing controlled noise to training data in order to simulate perturbations. The fraction of the data to be perturbed is determined by η . The noise to be added is generated with uniform random values scaled to the range [-2,2]. Table 12 shows that the SVQ-equiped model experiences a much lower rise in MSE and MAE relative to the model without SVQ, regardless of the proportion of injected noise. For instance, when the proportion of injected noise is 10%, the MSE of the model without SVQ rises by 25.4%, whereas the model with SVQ shows a modest increase of just 3.2%. These results confirm the effect of SVQ on helping the forecasting model handle noise better through vector quantization.

Noise proportion	MS	SE	M.	AE
rioise proportion	w/o SVQ	w SVQ	w/o SVQ	w SVQ
η=0	1.105	1.018	0.6567	0.6109
$\eta = 10\%$	1.386(+25.4%)	1.051(+3.2%)	0.7702(+17.3%)	0.6269(+2.6%
$\eta = 20\%$	1.554(+40.7%)	1.196(+17.5%)	0.8282(+26.1%)	0.6710(+9.89
$\eta = 30\%$	1.750(+58.4%)	1.255(+23.2%)	0.8821(+34.3%)	0.6953(+13.89
$\eta=40\%$	2.081(+88.3%)	1.568(+54.0%)	1.0031(+52.7%)	0.7973(+30.59
$\eta = 50\%$	2.646(+139.4%)	1.529(+50.2%)	1.1339(+72.7%)	0.7881(+29.09

Table 12: Noise injection analysis: The proportion of injected noise is indicated by η . We present MSE and MAE, and their percentage increase over baseline without artificial noise.

D.3 POSITION OF QUANTIZATION MODULE

To illustrate how the position of quantization module affects representation learning, we consider two designs as shown in Figure 8. In our main experiments, we adopt the first design where quantization is performed before the translator. In Section 5.4, we also investigate an alternative design where quantization is performed after the translator. Two designs only differ in the placement order of quantization module and translator module, while the other settings are kept the same.



Figure 8: Comparison of two quantization designs with different positions. Top: Quantization before translator. Bottom: Quantization after translator.

D.4 COMPREHENSIVE COMPARISON OF VQ METHODS USING THE SAME CODEBOOK SIZE

To make a fair comparison, we extend Table 4 by setting the codebook size to the same value (1024) 1063 for compared VQ methods. They are comprehensively evaluated from different aspects including 1064 downstream performance (prediction MSE), codebook usage (perplexity), and computational complexity (FLOPS, inference FPS, and training time per epoch). The quantitative results are shown in 1066 Table 13. The convergence performance is shown in Figure 9, where SVQ quickly converges to the 1067 lowest prediction error and satisfactory utilization of the codebook. Residual VQ with stochastic 1068 sampling has the highest codebook usage. However, its prediction MSE is worse than residual VQ 1069 without stochastic sampling. This demonstrates that forcibly improving codebook usage does not 1070 guarantee better downstream performance. SVQ generally outperforms the other VQ methods in 1071 computational efficiency, due to the simplified approximation described in Section 3.1.

Table 13: Quantitative comparison of VQ methods with the same codebook size (1024) on WeatherBench-S temperature dataset.

1075

1072

1058

1059

1061

1062

1039

1040

1076	Vector quantization method	Prediction MSE↓	Perplexity↑	FLOPS↓	Inference FPS \uparrow	Training time per epoch(min) \downarrow
1077	VQ	1.8544	51.95	7.207G	21.1	7.11
1079	Residual VQ	1.2131	142.47	8.616G	7.7	13.25
1070	Residual VQ (Stochastic)	1.8882	817.91	8.616G	8.1	17.27
1079	Grouped Residual VQ	1.1737	132.57	8.616G	4.8	19.98
	Multi-headed VQ	1.2113	16.36	8.717G	6.2	13.15
	SVQ (Frozen)	1.0393	335.72(0=3)/438.39(0=2)	8.037G	24.6	7.27
	SVQ (Learnable)	1.0403	$246.44(\theta=3)/331.41(\theta=2)$	8.037G	24.9	7.30



Figure 9: Prediction error and codebook usage of different VQ methods during the training process. All methods adopt the same codebook size (1024) and are early stopped with a patience of 10, trained on WeatherBench-S temperature dataset. For simplicity, the perplexity of SVQ is averaged on different θ .

1097 1098 D.5 Ablation of frozen module

The SVQ module consists of a two-layer MLP and a large codebook. The MLP can be seen as a projection from the input vector to the regression weights. In Section 5.5, we have already examined the impact of freezing the codebook on forecasting performance. To further investigate the effect of freezing the two-layer MLP, we conducted an ablation study in this section. The results, presented in Table 14, show that freezing the codebook only has a slight impact on forecasting performance, while freezing the MLP significantly impairs the performance. The MLP is essential for sparse regression and must be learned from the data, as it generates the weights needed to combine codes from the codebook.

1107

1108 1109 1110

1091

1092

1093

1094

1095 1096

1099

Table 14: Ablation of frozen modules on WeatherBench-S temperature dataset.

Metric		Frozer	n module	
	None (All learnable)	Codebook	Two-layer MLP projection	Both
MSE	1.018	1.023	1.060	1.093
MAE	0.6109	0.6131	0.6194	0.6387

D.6 EXPERIMENT STATISTICAL SIGNIFICANCE

To get more robust experimental results and evaluate the statistical significance, we rerun SimVP and SimVP+SVQ models five times under identical conditions. The results are presented without standard deviations in the main text due to space constraints. The results with standard deviations on the WeatherBench-S temperature, TaxiBJ, and WeatherBench-M datasets are reported in Tables 15, 16, and 17, respectively. The standard deviations of the SimVP+SVQ model are generally smaller than or comparable to those of the SimVP model.

1118

1126

1130

1127Table 15: Statistical significance of1128models on the WeatherBench-S tem-1129perature dataset.

Table 16: Statistical significance of models on the TaxiBJ dataset.

1131				Model	MSE↓	MAE↓	SSIM↑	PSNR↑
	Model	MSE↓	MAE↓	SimVP (w/o VO)	0.3246±0.0173	15.03 ± 0.26	0.9844 ± 0.0008	39.71±0.11
1132	SimVP (w/o VQ)	1.105 ± 0.043	$0.6567 {\pm} 0.0185$	SimVP+SVQ (Frozen)	0.3171 ± 0.0056	$14.68 {\pm} 0.03$	$0.9848 {\pm} 0.0001$	$39.83 {\pm} 0.01$
1100	SimVP+SVQ (Frozen)	1.023 ± 0.007	$0.6131 {\pm} 0.0050$	SimVP+SVQ (Learnable)	0.3191±0.0020	$14.64{\pm}0.02$	$0.9849 {\pm} 0.0002$	$39.86 {\pm} 0.02$
1155	SimVP+SVQ (Learnable)	1.018 ± 0.009	0.6109 ± 0.0064					

¹¹²⁵

1135

Variable	Tem	perature	Hur	nidity		Wind Co	omponent	Total Cl	oud Cover
Model	MSE↓	MAE↓	MSE↓	MAE↓		MSE↓	MAE↓	MSE↓	MAE↓
SimVP (w/o VQ) SimVP+SVQ (Frozen) SimVP+SVQ (Learnabl	$\begin{array}{c c} & 4.833 \pm 0.031 \\ & 4.427 \pm 0.017 \\ e & 4.433 \pm 0.021 \end{array}$	1.5246 ± 0.0077 1.4160 ± 0.0051 1.4164 ± 0.0054	340.06±0.44 360.15±0.55 360.53±0.81	12.738±0 12.445±0 12.449±0	030 24. 024 23. 013 23	535±0.076 915±0.098 908±0.062	3.4882±0.0086 3.4078±0.0033 3.4060±0.0029	25.332±0.052 24.968±0.126 24.983±0.115	3.4509±0.0 3.4117±0.0 3.4095±0.0
5	5) 1105±01021	1110120.0001	0000010001	12.119 ±0	010 201		5.1000±0.0025	21000201110	51107512010
D.7 BENCHN	IARK ON T	AXIBJ DAT	ASET						
Derteini			IDEI						
Tab	e 18: Perfo	rmance cor	nparison f	or SVQ	and b	aseline	models on	Tax1BJ.	
		Model		MSE↓	MAE↓	SSIM↑	PSNR↑		
	Con E3D-	vLSTMSHI et a	al. (2015) al. (2019a)	0.3358	15.32 14.98	0.9836	39.45 39.64		
	PhyD	NetGuen & The	ome (2020)	0.3622	15.53	0.9828	39.46		
	Pre	dNetLotter et a dRNNWang et a	1. (2017) al. (2017)	0.3516	15.91 15.31	0.9828 0.9838	39.29 39.51		
	M	IMWang et al.	(2019b) (2021)	$\frac{0.3110}{0.3268}$	14.96	0.9847	39.65 39.52		
	D	MVFNHu et al.	. (2023)	3.3954	45.52	0.8321	31.14		
	Predl PredF	RNN++Wang et RNN.V2Wang e	t al. (2018) et al. (2022)	0.3348	15.37 15.55	0.9834 0.9826	39.47 39.49		
	SimVI	TAUTan et al. (2	2023a)	0.3108	14.93	0.9848	39.74		
	Sinivi	imVP+SVQ (F	(2022) rozen	0.3240	<u>14.68</u>	0.9844 0.9848	<u>39.83</u>		
	Sin	1VP+SVQ (Lea Improveme	arnable) ent	0.3191 ↑ 1.7%	14.64 14.6%	0.9849 ↑0.1%	39.86 ↑0.4%		
				1 1 2 0 7 0	1=00.00	1002.00			
D.8 BENCHN	iark on W	/EATHER B	ENCH-HN	AV DAT	ASET				
Die Direin									
To evaluate the	performanc	e of SVO o	n high-din	nension	al data	. we con	nducted ad	ditional ext	berimen
utilizing the We	atherBench	ı dataset Ra	usp et al. (2	2020) ir	ı a Hig	h-dime	nsional Mu	ılti-Variabl	e (HMV
setting. This da	taset, relate	ed to real-w	orld weat	ther for	ecastir	ig, cons	sists of vari	ous meteo	rologic
variables that c	ontribute to	a total of 1	10 chann	els. The	ese inc	lude ter	mperature	at 2 m heig	t aboy
surface (t2m), v	fractional	agitude-dire	r(tcc) ho	0 m nei	gnt (u	10), acc	umulated r	1 vorticity	(ny)
Notably, severa	l variables a	are structure	ed across	multiple	e vertic	al laver	rs. For inst	ance, pv 50	(pv), ci
the potential vo	rticity at 50	hPa. The V	VeatherBe	nch-HN	IV dat	aset is s	imilar to W	leatherBen	ch-M b
includes signifi	cantly more	e channels	(increasin	ng from	4 to 1	10). It	is designed	l for multi	-variabl
six-hour interva	l forecasting	g, trained or	ı data fron	n 1980-1	2015, y	alidated	d on data fr	om 2016, a	ind teste
on data from 20)17-2018. C)wing to co	nstraints o	on page	length	1, we di	vided the p	erformanc	e metric
01 SIM VP and S	ov Q across	unese 110 () achieved	inannels ii	nto two e enhan	separa	the table t_{0}	7% on the	initial 55	es 19 an
while the impre	ovement or	the subsec	auent 55 a	channel	s was	6.1%	Consequer	ntly, the cu	mulativ
average improv			1			/0.		<i>j</i> , ine eu	
wordge miprov	ement acros	ss all 110 cl	hannels wa	as 8.9%	. Thes	e findin	igs undersc	ore that, do	espite th

Table 17: Statistical significance of models on the WeatherBench-M dataset.

complexities introduced by high-dimensional datasets, SVQ effectively adapts to and enhances the predictive capabilities of the backbone model.

1190 -	Channel	SimVP (w/o VQ)	SimVP+SVQ (Frozen)	SimVP+SVQ (Learnable)	Improvement
100	u10	2 1557	2 0849	2 0757	37%
192	v10	2.1557	2.0049	2.0757	2.4%
193	t2m	2.1013	2.1101	2.1117	11.0%
194	tisr	1 13E+05	32640	30756	72.7%
195	tee	0.22906	0.22028	0.22316	3.8%
196	tn	0.000147	9.96E-05	9.96E-05	32.1%
107	z 50	371 19	268.81	269.45	27.6%
197	z 100	333.31	243.12	245.07	27.1%
198	z 150	339.15	265.88	267.14	21.6%
199	z_200	376.12	303.11	302.84	19.5%
200	z_250	404.46	340.07	338.73	16.3%
201	z_300	404.8	347.96	347.64	14.1%
202	z_400	357.33	314.15	313.48	12.3%
202	z_500	309.93	272.9	270.52	12.7%
203	z_600	276.35	241.09	237.15	14.2%
204	z_700	252.58	218.67	215.86	14.5%
205	z_850	229.2	203.88	203.05	11.4%
206	z_925	229	208.52	206.46	9.8%
207	z_1000	242.2	220.18	219.18	9.5%
208	pv_50	3.97E-06	2.93E-06	2.92E-06	26.3%
200	pv_100	1.44E-06	1.32E-06	1.33E-06	7.8%
209	pv_150	9.72E-07	8.42E-07	8.44E-07	13.4%
210	pv_200	1.00E-06	8.64E-07	8.50E-07	15.0%
211	pv_250	1.05E-06	9.77E-07	9.73E-07	7.7%
212	pv_300	8.51E-07	7.88E-07	8.00E-07	7.4%
213	pv_400	3.61E-07	3.42E-07	3.41E-07	5.6%
21/	pv_500	2.35E-07	2.28E-07	2.28E-07	3.1%
217	pv_600	3.89E-07	3.31E-07	3.32E-07	14.8%
215	pv_700	6.97E-07	5.30E-07	5.27E-07	24.4%
216	pv_850	8.81E-07	7.70E-07	7.62E-07	13.6%
217	pv_925	1.00E-06	9.66E-07	9.58E-07	4.6%
218	pv_1000	1.40E-06	1.35E-06	1.34E-06	4.5%
219	r_50	1.9046	1.2558	1.391	34.1%
220	r_100	7.6728	6.0793	6.3087	20.8%
220	r_150	9.6361	8.4854	8.7663	11.9%
221	r_200	14.653	13.175	14.435	10.1%
222	r_250	19.433	18.728	18.964	3.6%
223	r_300	20.000	19.791	19.851	1.4%
224	r_400	19.78	19.375	19.330	2.1%
225	r_500	19.100	10.0/0	10.823	1./%
226	r 700	10.200	17.701 17.041	17.930	1.3%
207	r 850	17.210	17.041	17.033	-0.7%
227	r 925	10 086	9.613	0 6727	-0.770 4.702
228	r = 1000	7 5008	9.015 8 0027	7 2011	+ ./70 28 %
229	a 50	9 03F-08	7 40F-08	7.2711 7.55F-08	18 10%
230	q_00	9.05E-00	1.87F-07	1.87E-07	5 30%
231	q_100 q_150	1.97E-07	1.07E-07	1.07E-07	2.5%
	q_100	6.68F-06	6.48F-06	6 39 5-06	43%
232	q_200	2 38F-05	2.24F-05	2.26F-05	57%
233	q_200	5.61F-05	5 30F-05	5 40F-05	5.7%
234	q_300 q_400	0.000181	0.000174	0.000177	3.9%
235	q_==00	0.000375	0.000174	0.000362	3.4%
236	q_500 q_600	0.000592	0.000557	0.000502	6.6%
227	a 700	0.000826	0.000775	0.000784	6.2%
	1-,00	0.000020	•	0.000701	
230 239 -			Average improvement		11.7%

Table 19: MAE comparison for SimVP and SVQ on WeatherBench-HMV (The first 55 channels).

1240

4 - 5	Channel	SimVP (w/o VQ)	SimVP+SVQ (Frozen)	SimVP+SVQ (Learnable)	Improvement
6	a_850	0.001054	0.000989	0.000994	6.1%
-	q_925	0.0009	0.000757	0.00075	16.7%
1	a_1000	0.001024	0.000776	0.000767	25.1%
8	t_50	1.7426	1.2564	1.3594	27.9%
9	t_100	1.5365	1.2852	1.3544	16.4%
0	t_150	1.5825	1.3185	1.3377	16.7%
1	t_200	1.7535	1.6099	1.6243	8.2%
	t_250	1.7895	1.5499	1.5571	13.4%
2	t_300	1.5555	1.3903	1.3658	12.2%
3	t_400	1.6055	1.4029	1.3737	14.4%
4	t_500	1.6277	1.4652	1.4493	11.0%
5	t_600	1.693	1.5669	1.5636	7.6%
6	t_700	1.8961	1.7623	1.7412	8.2%
7	t_850	2.1494	2.0246	1.8345	14.6%
	t_925	2.2241	1.8933	1.8747	15.7%
8	t_1000	2.2335	1.891	1.834	17.9%
9	u_50	2.9882	2.5862	2.5981	13.5%
0	u_100	3.4492	3.3665	3.3759	2.4%
	u_150	3.9731	3.7005	3.7262	6.9%
J	u_200	4.8379	4.597	4.5959	5.0%
	u_250	5.725	5.5288	5.5195	3.6%
1	u_300	5.9494	5.7551	5.7492	3.4%
	u_400	5.1291	4.9461	4.9373	3.7%
	u_500	4.2227	4.0826	4.0643	3.8%
	u_600	3.7165	3.6228	3.5999	3.1%
	u_700	3.4591	3.3916	3.3759	2.4%
	u_850	3.0929	3.0253	3.0211	2.3%
	u_925	2.9/11	2.8879	2.8847	2.9%
	u_1000	2.3824	2.298	2.2888	3.9%
	V_50	2.4844	2.3521	2.346	5.6%
	V_100	3.0001	2.9852	2.983/	2.1%
	V_150	5.8109	3.0097	3.0702	5.1%
	V_200	4.7779	4.0044	4.3901	3.8%
	v_230 v_300	5.0021	5.3004	5 7501	2.9%
	v_300 v_400	5.0538	J.7084 4 9421	4 0351	2.4%
	v_ 1 00 v_500	4 1298	4.0279	4.0228	2.5%
	v_500 v_600	3 612	3 5264	3 5224	2.5%
	v 700	3 3112	3 2528	3 2485	1.9%
	v 850	3 0267	2 9877	2 9857	1.9%
	v 925	2.9499	2.8942	2.8931	1.9%
	v_1000	2.3969	2.3354	2.3381	2.6%
	vo.50	8.64E-06	8.50E-06	8.50E-06	1.6%
	vo_100	1.16E-05	1.15E-05	1.15E-05	0.5%
	vo_150	1.59E-05	1.58E-05	1.58E-05	0.9%
	vo_200	2.26E-05	2.24E-05	2.24E-05	0.8%
	vo_250	3.24E-05	3.22E-05	3.22E-05	0.6%
	vo_300	3.85E-05	3.82E-05	3.82E-05	0.9%
	vo_400	3.61E-05	3.58E-05	3.57E-05	1.0%
	vo_500	2.99E-05	2.96E-05	2.96E-05	0.8%
	vo_600	2.69E-05	2.67E-05	2.68E-05	0.8%
	vo_700	2.68E-05	2.65E-05	2.65E-05	1.1%
	vo_850	2.76E-05	2.76E-05	2.76E-05	0.2%
	vo_925	2.64E-05	2.63E-05	2.63E-05	0.2%
	1000	2.095.05	2 09E 05	2 08E 05	0.10%
	vo_1000	2.08E-05	2.08E-03	2.06E-05	0.1%

Table 20: MAE comparison for SimVP and SVQ on WeatherBench-HMV (The last 55 channels).

1296 E EFFECT OF SVQ ON LATENT REPRESENTATION

We investigated the impact of SVQ on the sparsity of regression weights in Figure 7. To further investigate its effect on the latent representation, we compared the distribution of batch tensors before and after applying SVQ. The tensors were transformed into normalized vectors, and their density distributions were estimated. As depicted in Figure 10, the representation after SVQ demonstrates a more compact distribution, indicating improved robustness to noise. These results further prove that SVQ can enhance forecasting performance by effectively handling noise in the data.

Figures 11, 12, 14, 15, and 13 present the comparison of latent feature maps before and after applying SVQ. These figures illustrate that the difference between foreground and background in the feature maps increases after SVQ. For example, in the KittiCaltech dataset, a clear distinction is observed between road conditions and sky (Figure 11). Similarly, in the WeatherBench-S temperature dataset, distinctive regions are identified between high and low latitudes (Figure 12). These findings suggest that SVQ helps in enhancing the discriminative power of the latent representations, which in turn contributes to improved downstream forecasting performance.



Figure 10: Distribution of latent vector on WeatherBench-S temperature dataset.



Figure 11: Latent feature map on the KittiCaltech dataset: (a) feature map before SVQ, and (b) feature map after SVQ.



after SVQ.

			(a) Beto	re SVQ							(D) AII	er svQ			
Channel 1	Channel 2	Channel 3	Channel 4	Channel 5	Channel 6	Channel 7	Channel 8	Channel 1	Channel 2	Channel 3	Channel 4	Channel 5	Channel 6	Channel 7	Channel 8
Channel 9	Channel 10	Channel 11	Channel 12	Channel 13	Channel 14	Channel 15	Channel 16	Channel 9	Channel 10	Channel 11	Channel 12	Channel 13	Channel 14	Channel 15	Channel 16
	1.000		10	Dares.	16	- Aller					26 Km	() TO CO	15 A.#		
Channel 1	7 Channel 18	Channel 19	Channel 20	Channel 21	Channel 22	Channel 23	Channel 24	Channel 17	Channel 18	Channel 19	Channel 20	Channel 21	Channel 22	Channel 23	Channel 24
Channel 2	Channel 26	Channel 27	Channel 28	Channel 29	Channel 30	Channel 31	Channel 32	Channel 25	Channel 26	Channel 27	Channel 28	Channel 29	Channel 30	Channel 31	Channel 32
1000	ise 20	10.00	10	101		A BOARD	In line	A STATE	in the	<u>lieka</u>			16	ign -	Jan H
Channel 3	3 Channel 34	Channel 35	Channel 36	Channel 37	Channel 38	Channel 39	Channel 40	Channel 33	Channel 34	Channel 35	Channel 36	Channel 37	Channel 38	Channel 39	Channel 40
Channel 4	Channel 43	Channel 43	Channel 44	Channel 45	Channel 46	Channel 47	Channel 48	Channel 41	Channel 42	Channel 43	Channel 44	Channel 45	Channel 46	Channel 47	Channel 48
10	10.00	10	16000			19.950		ie i =	Jesse -	16 Ber	10000		Cistored		2000
Channel 4	Channel 50	Channel 51	Channel 52	Channel 53	Channel 54	Channel 55	Channel 56	Channel 49	Channel 50	Channel 51	Channel 52	Channel 53	Channel 54	Channel 55	Channel 56
Channel 5	7 Channel 58	Channel 59	Channel 60	Channel 61	Channel 62	Channel 63	Channel 64	Channel 57	Channel 58	Channel 59	Channel 60	Channel 61	Channel 62	Channel 63	Channel 64
	10 84	10C	1000	1560	56 C			(ene	1000	1670		1	15.000	Ser.	Sector.
Figure	e 15: 1	Latent	featur	re maj	p on th	he Hu	man3.	6M dat	taset:	(a) fe	ature	map b	oefore	SVQ	, and (b
ieatur	e map	after S	SVQ.												
					TATT	VFR	COLL	ГS							
F A	ADDI'	ΓΙΟΝ	AL O	UALI	IAH	Y LL IN.	ESUL.	10							
F A	ADDI'	ΓION	AL Q	UALI	IAII	VL R	ESUL.	10							
F A	ADDI' Fore	ΓΙΟΝ Casti	AL Q	UALI .RORS	ON W	VE K	HERBE	NCH A	ND TA	axiBJ	DATA	SETS			
F A	ADDI'	ΓΙΟΝ Casti	AL Q	UALI .RORS	ON W	VE K.	IERBE	NCH A	nd Ta	axiBJ	DATA	SETS.			
F A	DDI' Fore	ΓΙΟΝ CASTI	AL Q	UALI RORS	ON W	VEATH	HERBE	NCH A	nd Ta	AXIBJ	DATA	SETS			
F A	DDI' FORE 2 Frames 12 Frames	ΓΙΟΝ CASTI	AL Q	UALI RORS	ON W	VEATH	IERBE	NCH A	ND TA	AXIBJ	DATA	SETS			
F A	DDI' Fore 2 Frames 12 Frames	TION CASTI	AL Q	UALI	ON W	VEATH	HERBE	NCH A	ND T	AXIBJ	DATA	SETS	La Participa de		
F.1	DDI' Fore 2 Frames 12 Frames	TION CASTI	AL Q	UALI	ON W	VEATH	IERBE	NCH A	nd Ta	AXIBJ	DATA	SETS			
F.1	DDI' FORE 2 Frames 12 Frames STM Erre	ΓΙΟΝ CASTI	AL Q	UALI	ON W	VEATH	IERBE	NCH A	nd Ta	AXIBJ	DATA	SETS			
F A	2 Frames 12 Frames TM Erro	ΓΙΟΝ CASTI	AL Q	UALI	ON W	VEATH	IERBE	NCH A	nd Ta	AXIBJ	DATA	SETS			
F A F. 1 Input (1 Target (ConvLS E3D-LS PredRN	DDI' FORE 2 Frames 12 Frames STM Error STM Error	TION CASTI	AL Q	UALI	ON W	VEATH	IERBE	NCH A	nd Ta	AXIBJ	DATA	SETS			
F A F. 1 Input (1) Target (ConvLS E3D-LS PredRN	DDI' FORE 2 Frames 12 Frames STM Error STM Error	ΓΙΟΝ CASTI	AL Q	UALI	ON W	VEATH	IERBE	NCH A	nd Ta	AXIBJ	DATA	SETS			
F A F. 1 Input (1 Target (ConvLS E3D-LS PredRN MAU E	DDI' FORE 2 Frames 12 Frames 12 Frames TM Error STM Error IN Error	ΓΙΟΝ CASTI	AL Q	UALI	ON W	VEATH	IERBE	NCH A	nd T	AXIBJ	DATA	SETS			
F A F. 1 Input (1 Target (ConvLS E3D-LS PredRN MAU E PredRN	2 Frames 2 Frames 12 Frames STM Error STM Error IN Error IN Error	ΓΙΟΝ CASTI	AL Q	UALI	ON W	VEATH	IERBE	NCH A	ND T	AXIBJ	DATA	SETS			
F A F.1 Input (1 Target (ConvLS E3D-LS PredRN MAU E PredRN	DDI' FORE 2 Frames 12 Frame STM Error STM Error IN Error IN Error	ΓΙΟΝ CASTI s) r r	AL Q	UALI RORS	ON W	VEATH	IERBE	NCH A	nd Ta	AXIBJ	DATA	SETS			
F A F. 1 Input (1 Target (ConvLS E3D-LS PredRN MAU E PredRN PredRN	FORE 2 Frames 12 Frames 3TM Error 3TM Error	TION CASTI () () () () () () () () () () () () ()	AL Q	UALI	ON W	VEATH	IERBE	NCH A	nd Ta	AXIBJ	DATA	SETS			
F A F. 1 Input (1 Target (ConvLS E3D-LS PredRN MAU E PredRN PredRN SimVP	FORE 2 Frame 2 Frame 3 TM Error 3 TM Error	TION CASTI	AL Q	UALI	ON W	VEATH	IERBE	NCH A	ND T	AXIBJ	DATA	SETS			
F A F. 1 Input (1 Target (ConvLS E3D-LS PredRN MAU E PredRN MAU E	FORE 2 Frames 12 Frames 12 Frames STM Error STM Error IN Error IN Error IN V2 Error IN V2 Error IN ++ Error Error	ΓΙΟΝ CASTI () () (s) (s) (s) (s) (s) (s) (s) (s) (AL Q	UALI	ON W	VEATH	IERBE	NCH A	nd Ta	AXIBJ	DATA	SETS			
F A F. 1 Input (1 Target (ConvLS E3D-LS PredRN MAU E PredRN PredRN SimVP Ours E	FORE 2 Frames 2 Frames 12 Frames 3TM Error 3TM Error	TION CASTI () () (s) (s) (s) (s) (s) (s) (s) (s) (AL Q	UALI	ON W	VEATH	IERBE	NCH A	ND TA	AXIBJ	DATA	SETS			

Figure 16: The qualitative forecasting errors on WeatherBench-S temperature dataset.



