

3D Motion Perception of Binocular Vision Target with PID-CNN

Shi Jiazhao*, Pan Pan, Shi Haotian

Xi'an MTC Institute

*To whom correspondence should be addressed; E-mail: shijiazhao123@stu.xjtu.edu.cn

Abstract—This article trained a network for perceiving three-dimensional motion information of binocular vision target, which can provide real-time three-dimensional coordinate, velocity, and acceleration, and has a basic spatiotemporal perception capability. Understood the ability of neural networks to fit nonlinear problems from the perspective of PID. Considered a single-layer neural network as using a second-order difference equation and a nonlinearity to describe a local problem. Multilayer networks gradually transform the raw representation to the desired representation through multiple such combinations. Analysed some reference principles for designing neural networks. Designed a relatively small PID convolutional neural network, with a total of 17 layers and 413 thousand parameters. Implemented a simple but practical feature reuse method by concatenation and pooling. The network was trained and tested using the simulated randomly moving ball datasets, and the experimental results showed that the prediction accuracy was close to the upper limit that the input image resolution can represent. Analysed the experimental results and errors, as well as the existing shortcomings and possible directions for improvement. Finally, discussed the advantages of high-dimensional convolution in improving computational efficiency and feature space utilization. As well as the potential advantages of using PID information to implement memory and attention mechanisms.

Index Terms—3D Motion Perception; Binocular Vision; PID-CNN; High-dimensional Convolution; Memory and Attention Mechanisms

I. INTRODUCTION

Coordinate, velocity, and acceleration are very useful for many control problems^[1], but in actual measurement, often face constraints from measurement tools and timeliness. Directly measuring coordinate using a measuring ruler and then calculating velocity and acceleration has low requirement for measuring tool. But it is difficult to measure targets moving at high speed, and has poor timeliness. Indirect measurement can be classified according to the measurement principle as echo measurement, visual measurement, and holographic measurement^[2, 3, 4].

Echo measurement is the process of calculating distance using the echo time of the measurement signal, and calculating three-dimensional coordinate of the target based on the azimuth. Then use the Doppler effect or directly use coordinate to calculate velocity and acceleration. Radar, sonar, laser, and infrared measurements all belong to this category. This type of measurement requires actively emitting measurement signal, such as electromagnetic waves or sound waves, and scanning

different spatial orientations through rotation, phase control, or other methods. The constraints of its measurement tools mainly come from the signal emitting and reception devices, and the constraints of timeliness mainly come from the time of signal emitting and reception, the speed of signal propagation, and the timeliness of data post-processing algorithms. Visual measurement restores three-dimensional information through the parallax of different perspectives. A single perspective image is two-dimensional projection of three-dimensional space, losing depth dimension information. Two images from different perspectives provide four dimensional constraints, allowing for the restoration of three-dimensional information. Comparing the above two measurement types, echo measurement directly measures information of the depth dimension, and then supplements information of the other two dimensions by adjusting the emitting azimuth. Visual measurement directly measures information of the height and width dimensions, and then supplements information of the depth dimension through parallax calculation. Holographic measurement can be seen as a combination of these two methods, simultaneously recording the amplitude and phase of the interference waves, which is equivalent to recording information of a third dimension in addition to height and width. However, it does not record colour information and has higher constraints on the measurement tools.

Comparing the auditory and visual senses of animals, it can be considered that hearing passively perceives three-dimensional information in the first way, and vision perceives three-dimensional information in the second way. Only a few species such as bats and dolphins have the ability to actively perceive in the first way. Obviously, vision provides more information, but it relies on external lighting conditions, and when lighting is insufficient, hearing plays an important supplementary role. In most industrial applications, lighting is a controllable condition. Meanwhile, due to the advancement of camera equipment, the constraints of visual measurement tools compared to those of echo measurement tools has been greatly reduced. The constraints of visual measurement timeliness mainly come from the vision signal acquisition time, light propagation time, and the timeliness of data post-processing algorithms. In the first two time constraints, the echo and visual measurement methods are basically equivalent, so the data post-processing algorithm becomes the key factor affecting their timeliness. Most data post-processing algorithms for echo

measurement are based on Fourier transform^[5, 6, 7], which can achieve real-time high-speed target measurement.

The data post-processing algorithm for visual measurement is mainly based on the principle of triangulation^[8]. Single perspective observation can limit the target to a line of sight emitted from the observation point, and another different perspective can provide another line of sight. By calculating the intersection of these two lines, the three-dimensional coordinate of the target can be obtained. So it is necessary to first calibrate the geometric relationship between two observation points, measure the azimuth angles of the two lines of sight, and then obtain the coordinate of the target in the world coordinate system through the principle of triangulation and coordinate transformation. This process can be equivalent to a multi-step matrix transformation. It is not difficult to obtain the parameters of the transformation matrix after calibration, but calibration is often tedious and time-consuming, which greatly affects the timeliness of visual measurement. Matrix transformation is exactly what neural networks excel at^[9, 10, 11], and we hope to replace the tedious and time-consuming calibration process with neural networks.

For this, we trained a PID convolutional neural network for perceiving three-dimensional motion information of a binocular vision target. The perception target was a 3D randomly moving ball generated by simulation. First, we analysed some useful reference principles for designing neural networks' architecture. Understood the ability of neural networks to fit nonlinear problems from the perspective of PID. Implemented a simple but practical feature reuse method. The model is relatively small, and converged quickly during training. The experimental results shown that the network can predict the coordinate, velocity and acceleration of the target in real time, with an accuracy close to the upper limit that the input image resolution can represent. Analysed the experimental results and errors, as well as the limitations and possible improvements. Finally, we discussed the advantages of high-dimensional convolution in improving computational efficiency and feature space utilization. And explained the advantages of using PID information to implement attention and memory mechanisms.

II. NETWORK DESIGN ANALYSIS

Most higher animals in nature have evolved binocular vision, indicating the advantage of binocular vision over monocular vision in providing three-dimensional information, as well as the efficiency advantage over multi-ocular vision. Image records the frequency and amplitude of the reflected light from objects, corresponding to the visual colour and intensity. A colour image with a height and width of 224 pixels can be represented as a tensor of $3*224*224$. The dimension of tensor often causes confusion in different fields. To avoid this confusion, this article refers to the tensor above as a three-dimensional tensor, with the first dimension having 3 degrees of freedom, and the second and third dimensions having 224 degrees of freedom.

The usual convolutional transformation for images is to convolve the height and width dimensions with a kernel of size 3, without making a distinction between the colour and feature

dimensions^[12, 13, 14], performing fully connected transformations on those dimensions. For example, considering a $3*224*224$ feature map been transformed to a $24*224*224$ feature map through 24 convolutional kernels, usually, the first dimension of resulting feature map is understood as 24 feature channels, replacing the colour channels and adding degrees of freedom. If distinguished, it can be seen as obtaining a $8*3*224*224$ feature map, still retaining the colour dimension and forming a new feature dimension with 8 degrees of freedom. When performing subsequent convolution transformations, it simultaneously observes three adjacent pixels in height and width dimensions, three colour channels, and eight feature spaces. Subsequent transformations can also be seen as forming new feature dimensions, rather than adding degrees to the existing feature dimension. There is no problem not making such a distinction when using fully connected transformations on those dimensions. But if higher dimensional convolutions are used^[15], this distinction is necessary.

Many experiments have shown that convolutional kernel of size 3 has greater advantages compared to other sizes^[16, 17, 18]. We are trying to understand it from the perspective of PID^[19]. In the field of control, PID represents through proportional, integral, and derivative transformations of input signal to obtain control signal. Proportion represents the response that the system should make to the current input. Integral represents the response that the system should make to the accumulation of input. Derivative represents the response that the system should make to the changing rate of input. For one dimension situation, a convolutional kernel of size 3 can be regarded as a PID signal extractor. The convolutional kernel C_p of $[0, 1, 0]$ can extract proportional information. Average pooling is equivalent to the convolutional kernel C_i of $[1/3, 1/3, 1/3]$, which can extract integral information. The convolutional kernel C_d of $[-1, 0, 1]$ can extract first-order difference information. By combining different PID coefficients k_p, k_i, k_d , a convolutional kernel: $C = k_p * C_p + k_i * C_i + k_d * C_d$ can be obtained, which can extract information of interest. When $k_p = -3, k_i = 3, k_d = 0$, a convolutional kernel C_{d2} of $[1, -2, 1]$ can be obtained, which can extract second-order difference information.

Many fields use second-order differential information, and rarely use higher-order differential information. In the field of control, 2nd order controllers are usually used, and higher order controllers often have poorer stability. In mathematical modeling, combinations of second-order partial differential equations and nonlinearities are commonly used, and higher-order partial differential equations are rarely used. For dynamic problems, usually use up to second-order differential, i.e., acceleration, and pay less attention to higher-order information. These indicate that second-order differential information has significant value for many problems, and higher-order information may be impractical due to its complexity. Therefore, a convolutional kernel of size 3 is a suitable minimum choice for fitting most problems.

For a neural network with feature map of 1 dimension, let $f_L(x)$ denote the activation at position x in the L -th layer, and define discrete its first-, second-order difference and integral as:

$$\begin{cases} f_L'(x) = f_L(x+1) - f_L(x-1) \\ f_L''(x) = f_L(x+1) - 2f_L(x) + f_L(x-1) \\ \int_i^j f_L(x)dx = f_L(i) + \dots + f_L(j) \end{cases} \quad (1)$$

The transformation of a single-layer convolution and nonlinearity can be described as:

$$\begin{cases} k_p f_L(x) + \frac{1}{3} k_i \int_{x-1}^{x+1} f_L(x)dx + k_d f_L'(x) = z_L(x) \\ f_{L+1}(x) = g(z_L(x)) \end{cases} \quad (2)$$

where $z_L(x)$ is the intermediate variable, and $g(\cdot)$ is the nonlinear transformation. If C_i , C_d , and C_{d2} are used as base vectors, another form can be obtained:

$$\begin{cases} k_{p2} f_L(x) + k_{d2} f_L'(x) + k_s f_L''(x) = z_L(x) \\ f_{L+1}(x) = g(z_L(x)) \end{cases} \quad (3)$$

where:

$$\begin{cases} k_{p2} = k_p + k_i \\ k_{d2} = k_d \\ k_s = \frac{1}{3} k_i \end{cases} \quad (4)$$

From this perspective, it can be considered that a single-layer transformation is equivalent to using a second-order difference equation and a nonlinearity to describe a local problem. Multilayer transformations are equivalent to fitting complex nonlinear problems through multiple such combinations, transforming the input raw representation to the desired representation gradually.

After convolutional linear transformation, it is necessary to undergo a nonlinear transformation to obtain a new feature representation^[20, 21, 22]. Therefore, how to evaluate the quality of a representation is of great significance for designing the architecture of a network. Although it is difficult to define what is a good representation, there are some principles that are effective in most cases. A good representation is related to specific weights, has lower information entropy, and occupies less space. It retains the information content of interest and discards redundant information. Taking the classification problem as an example, original cat image occupies the most space, and has the highest information entropy. The label of "cat" occupies the least space and has the lowest information entropy. It is a good representation for those who have established the concept of "cat", and for those who only speak Chinese, "猫" is the corresponding good representation.

Figure 1 illustrates that through linear and nonlinear transformations, representing the information of interest more explicitly in a 2D space. To form a closed feature subspace in 2D space, three segmentation lines are required. Therefore, this transformation is unable to extract information from a closed subspace. There are two approaches, one is to increase the degrees after transformation, and the other is through multi-step transformations.

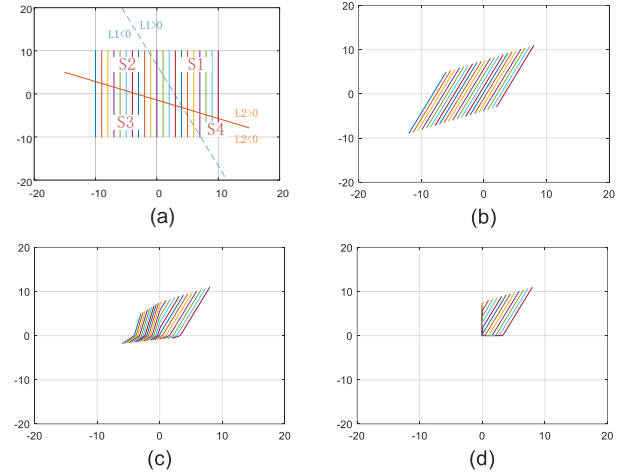


Fig. 1. Through linear and nonlinear transformations, representing the information of interest more explicitly. (a) is original representation ($L_1: w_{11}*x_1+w_{12}*x_2+b_1=0$, $L_2: w_{21}*x_1+w_{22}*x_2+b_2=0$, the parameters come from the linear transformation). (b) is the representation after PRelu, $\alpha=0.2$, it compresses the negative space. (d) is the representation after Relu, mapping the negative space to the origin and axes. Both (c) and (d) can represent the S1 subspace in (a) more explicitly.

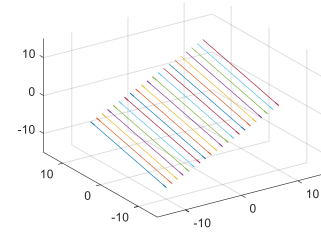


Fig. 2. Increase the degrees of freedom after transformation.

Figure 2 shows according to the first approach, mapping 2D space to 3D space. After mapping, the feature map plane itself limits one degree of freedom, so the number of planes required to enclose a closed feature map in 3D space is still 3. But, as the segmentation changes from lines to planes, the required parameters for subsequent transformations will increase. Increasing the degrees after transformation does not increase information content or reduce information entropy, nor does it make the segmentation of feature space easier. Therefore, we prefer to use multi-step transformations with equal degrees.

Using transformations with equal degrees, we hope not to reduce information content, that is, to be able to restore input information. The condition is that the linear transformation is full rank, and the nonlinearity is invertible. Comparing two types of nonlinearity, Relu can more explicitly represent the information of interest, but PRelu has other advantages. In addition to propagating gradients in negative space, PRelu is also a invertible mapping, so it does not reduce information content. Besides, for Figure 1. (c), when the positive space of the subsequent transformation is inversely mapped in the third quadrant in (c), the subsequent transformation can extract the information that was ignored by the previous transformation. This relationship is similar to that when the former neuron is in

suppression, the latter neuron is in activating, which may lead to better representation.

Extending the above analysis to cases with more degrees of freedom, it can be considered that linear transformation provides segmentation hyperplanes, and nonlinearity represents the subspace of interest more explicitly.

After representing the information of interest more explicitly, it need to use pooling to reduce the size of feature map, discard redundant information, i.e., reduce the information entropy. As mentioned earlier, we consider average pooling as an integral kernel. Subsequent experimental result has shown that average pooling converges faster on our problem. Figure 3 show a intuitive comparison, we think average pooling retains more information content.



(a) original image



(b) average pooling



(c) max pooling

Fig. 3. After four pooling operations and enlarge to original size, kernel size and stride are both 2.

The BN layer with parameters can move the feature space as a whole towards the desired subspace direction in backpropagation^[23], so its combination with linear transformation can accelerate convergence. We think it is more suitable to place it after the linear transformation, and before the nonlinearity represents information of interest explicitly.

Based on the above analysis, we have designed a PID convolutional neural network for perceiving three-dimensional motion information of binocular vision target.

III. NETWORK ARCHITECTURE

The building block of the network is shown in Figure 4. On the right is the reference feature map size. The input is a four-dimensional tensor, with the first dimension representing 2 perspectives, the second dimension representing time, and the last two dimensions representing the height and width. We did not use colour images because we think that colour plays a secondary role in motion perception. After two rounds of Conv-BN-PReLU transformations, information from a closed subspace can be extracted. The resulting feature map and input feature map are concatenating on the first dimension. Then use an average pooling with kernel size of 2 for height and width dimensions. Performing fully connected transformation on the first dimension, and using convolutional transformation with kernel size of 1 on the second dimension.

We used 7 building blocks to transform the input from $2*3*256*256$ to $256*3*2*2$. Due to the feature reuse, half of a block's input comes directly from the pooling result of the

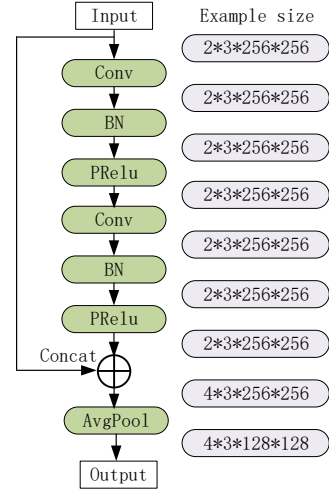


Fig. 4. Building Block

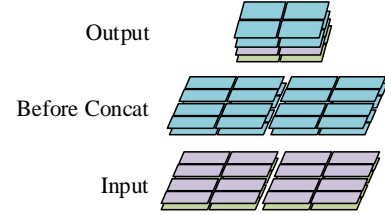


Fig. 5. Illustration of Feature Reuse

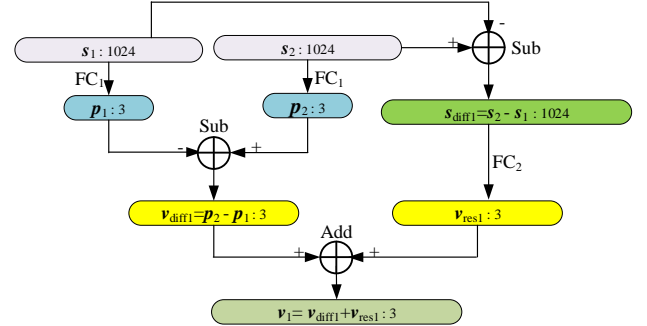


Fig. 6. Residual process for calculating velocity. Using the fully connected transformation FC_1 to obtain coordinate vectors p_1, p_2 . Calculating the difference between p_2 and p_1 , obtained $v_{diff1} = p_2 - p_1$. Calculating the difference between s_2 and s_1 , obtained $s_{diff1} = s_2 - s_1$. Then use another fully connected transformation FC_2 to obtain the velocity residual vector v_{res1} . The final velocity result is $v_1 = v_{diff1} + v_{res1}$. The number after the colon represents the vector dimension.

previous block's feature map before concatenation, and the other half comes from the pooling results of all other previous blocks', that is, 1/4 from the previous second block's, and so on. Figure 5 briefly illustrates the characteristic of this feature reuse. For a $2*2$ feature map in the last building block's output, the information of its first element comes from the top left quarter information of all previous blocks' output. This feature reuse method can bring significant benefits to both forward and backward propagation. During forward propagation, each block sees the pooling results of all previous blocks. During backpropagation, gradients can be directly propagated to the first block.

For the obtained $256*3*2*2$ feature map, we split and flatten it in the second dimension to obtain three 1024

dimension vectors: s_1 , s_2 , and s_3 , from which we calculate coordinate, velocity, and acceleration. Figure 6 shows a residual process of calculating velocity. This is done because, if the positional coordinates have been trained well, we hope to utilize this information, and also use the velocity loss to provide gradient feedback, without excessively changing the well trained weights. We also applied similar process to acceleration.

Overall, the network consists of 7 building blocks and 3 fully connected layers. Each building block contains two convolutional layers. Therefore the total number of layers is 17, and the total parameters is 413 thousand.

IV. DATASET AND TRAINING PROCESS

A. Dataset

The data we used is binocular vision images of a simulated randomly moving ball, as shown in Figure 7. The two observation points are located in the $(-4, -5, 5)$ and $(-5, -4, 5)$ directions of the projection center, the line of sight angle is 9.9866° . The coordinate axis range is $(-50, 50)$, and the diameter of the ball is 10. To prevent the ball from exceeding the field of view, the coordinate is uniformly randomly generated between $(-45, 45)$.

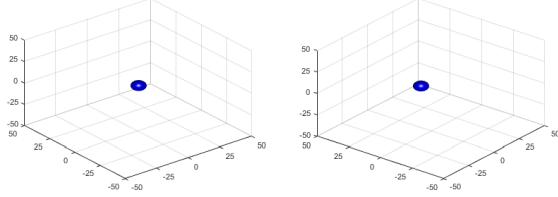


Fig. 7. Left perspective and right perspective, coordinate is (18.842835, 22.921801, -20.157743).

The train dataset consists of 11152 binocular images, the validation and test datasets consist of 1024 binocular images respectively. We take the time of one frame as one time unit. For 3 frames of input, the first frame is taken according to the index order in dataset, and the last two frames are taken randomly. So the range of velocity is $(-90, 90)$, and the range of acceleration is $(-180, 180)$. The training number of each epoch is the number of binocular images in the dataset.

When representing the coordinate range of $(-50, 50)$ in a 256×256 pixels image, the coordinate resolution is $100/256 = 0.3906$, which means that the ball must move at least 0.3906 to produce a pixel difference in the image. Considering the presence of boundaries and the image represents 3D space, the actual coordinate resolution that a single image can represent should be even lower.

For the input, we used the B channel in the RGB image and subtracted the background from each perspective, so the network does not see the coordinate axes, grids, and numbers in Figure 7. For the output, we normalized using the theoretical mean and variance of the coordinate distribution, although the theoretical variances of velocity and acceleration are different. No other data processing was done.

B Training Process

The training process has three stages. First, train a network

that only measures coordinate with a single frame input. Then use the former weights to train a network with a input of 2 frames, measures 2 coordinates and 1 velocity. Finally, train the complete network with input of 3 frames, which output is an 18 dimensional vector, representing 3 coordinates, 2 velocities and 1 acceleration vectors in sequence.

We use the MSE loss and Adam solver^[24], with a batch size of 32. When training the coordinate measurement network, the initial learning rate was $1e-3$, with an attenuation of 0.1 every ten epochs. After every four rounds of attenuations, learning rate was amplified $1e3$ times. Total epochs for training coordinate measurement network is 120. When training the velocity and acceleration measurement networks, the initial learning rate is $1e-6$. Other hyperparameter are same as the coordinate measurement network.

V. EXPERIMENTAL RESULTS AND ANALYSIS

The experimental results on the test dataset are shown in Table 1. The standard deviations of coordinate, velocity and acceleration prediction are 0.181698, 0.245950 and 0.426908, respectively. That means the network's prediction of coordinate: $\hat{p} \sim N(p^*, 0.181698^2)$. The coordinate accuracy is higher than the coordinate resolution we have estimated, which may be due to the higher accuracy brought by the four dimensions' constraints of binocular vision, and the additional accuracy provided by the transition area on the boundary of the ball. We think the accuracy has approached the upper limit that the input image resolution can represent. The maximum error of positional coordinate prediction is 0.700486, corresponding input images are shown in Figure 8. The maximum error of velocity and acceleration are 0.983588 and 1.708877.

TABLE 1
EXPERIMENTAL RESULTS OF TEST DATASET

	standard deviation	maximum error
coordinate	0.181698	0.700486
X axis	0.181393	0.640285
Y axis	0.176532	0.685821
Z axis	0.187017	0.700486
velocity	0.245950	0.983588
X axis	0.240515	0.809299
Y axis	0.248256	0.983588
Z axis	0.248990	0.834841
acceleration	0.426908	1.708877
X axis	0.418635	1.708877
Y axis	0.424857	1.532333
Z axis	0.437027	1.434120

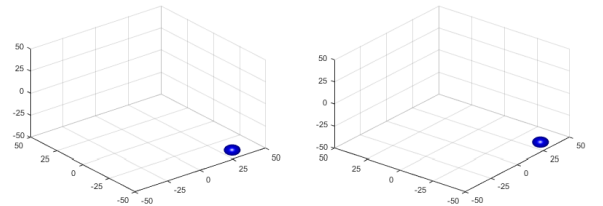


Fig. 8. Images with maximum coordinate prediction error. The real coordinate is (28.389164, -44.877795, -44.721836), and the predicted is (27.893255, -44.69415, -44.02135).

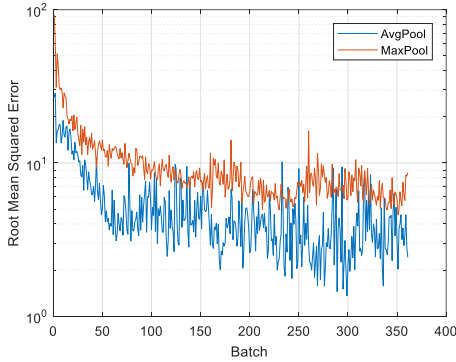
The errors after each stage are shown in Table 2. After 80th epoch of each stage, no obvious accuracy improvement was observed. It shown that after providing feedback on velocity loss and acceleration loss, both coordinate and velocity accuracy had been improved.

TABLE 2
ERRORS AFTER EACH STAGE

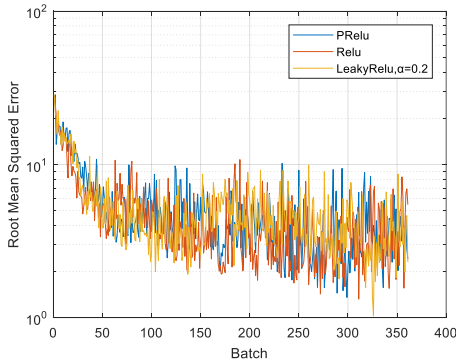
	standard deviation		
	coordinate	velocity	acceleration
First stage	0.251755	/	/
Second stage	0.202216	0.284662	/
Third stage	0.181698	0.245950	0.426908

In the third stage network, if don't use residual process, directly use coordinate to calculate velocity and acceleration, the standard deviation of velocity and acceleration are 0.252794 and 0.440178. The use of residual process reduced these standard deviations by 2.71% and 3.01% respectively.

Figure 9 (a) compares the convergence speed of average pooling and maximum pooling in the first epoch of training coordinate measurement network. It shows that average pooling is converging faster overall than maximum pooling. Figure 9 (b) compares the convergence speeds of different nonlinearities, no obvious differences were observed. We think this is due to that our problem does not have high requirement on nonlinearity.



(a) Comparison of different poolings' convergence speed



(b) Comparison of different nonlinearities' convergence speed

Fig. 9. Comparison of convergence speed

Performing 1024 tests on the test dataset consumes 4.142 s. The time for one measurement is 4.05 ms, and 247 measurements can be done per second. It can be considered that the model has achieved real-time measurement capability.

VI. CONCLUSION AND DISCUSSION

This article analysed some reference principles for designing neural networks used for perceiving three-dimensional motion, and based on this, designed a PID convolutional neural network. On the simulated dataset of randomly moving balls, the measurement accuracy is close to the upper limit that the input image resolution can represent. It can be considered has achieved real-time perception of three-dimensional motion information. However, there are still shortcomings, and improvement can be done as follows:

- 1) This article used simulated images of a single target from fixed perspectives. We hope to test on the real physical dataset and measure targets of different shapes, as well as measure relative coordinate from variable perspectives in the future.
- 2) The main process of triangulation is linear transformation, that's the main reason why we didn't observe the convergence speed differences of different nonlinearities. We hope to test the model on more complex nonlinear problems in the future.

Despite the above shortcomings, establishing a basic spatiotemporal perception capability for neural networks is still an important step.

Besides, in network design analysis section, we have mentioned higher dimensional convolution. We think it can play an important role in improving computational efficiency and effectively utilizing feature space. And we also explained the ability of neural networks to fit complex problems from the perspective of PID. But these have not been fully validated in this article yet.

Currently, most deep learning frameworks do not support convolutions higher than 3 dimensions. Our model uses fully connected transformation for the feature dimension, and the computational complexity of kernel grows exponentially by 4 times (ignoring bias here). If a high-dimensional convolution with kernel size of 3 is used, the growth rate can be reduced to 3, which has a significant impact on deeper networks. And an n-dimensional convolution with kernel size of 3 can simultaneously observe 3^n adjacent elements. These elements distribute closely in different orthogonal dimensions, and each element has more "neighbors", which allows the feature space to be more efficiently utilized.

Understanding the ability of neural networks from the perspective of PID, may play a greater role in the memory and attention mechanisms^[25, 26, 27]. Integral means historical information, difference means changing information. In dot product attention, assigning larger attentional weights to vectors have larger dot products. Two vectors which have smaller difference, have smaller Euclidean distances and larger dot product. By using difference information to perform weighted integration on proportional information, can achieve more sophisticated representation. Integration at different time scales can represent memory at different time scales. Integration is addition, difference operation is subtraction, and their computational complexity is obviously lower. We hope to do more in-depth research on the above aspects in the future.

REFERENCES

- [1] Jeon S, Tomizuka M. Benefits of acceleration measurement in velocity estimation and motion control. *Control Engineering Practice*. 2007 Mar 1;15(3):325-32.
- [2] Kuo J, Von Ramm OT. Three-dimensional motion measurements using feature tracking. *IEEE transactions on ultrasonics, ferroelectrics, and frequency control*. 2008 Apr 25;55(4):800-10.
- [3] Lin X, Wang J, Lin C. Research on 3D reconstruction in binocular stereo vision based on feature point matching method. In *2020 IEEE 3rd International Conference on Information Systems and Computer Aided Education (ICISCAE) 2020 Sep 27* (pp. 551-556). IEEE.
- [4] Tian X, Liu R, Wang Z, et al. High quality 3D reconstruction based on fusion of polarization imaging and binocular stereo vision. *Information Fusion*. 2022 Jan 1;77:19-28.
- [5] Xu J, Xia XG, Peng SB, et al. Radar maneuvering target motion estimation based on generalized Radon-Fourier transform. *IEEE Transactions on Signal Processing*. 2012 Sep 3;60(12):6190-201.
- [6] Chen VC, Qian S. Time-frequency transform vs. Fourier transform for radar imaging. In *Proceedings of Third International Symposium on Time-Frequency and Time-Scale Analysis (TFTS-96) 1996 Jun 18* (pp. 389-392). IEEE.
- [7] Su X, Chen W. Fourier transform profilometry:: a review. *Optics and lasers in Engineering*. 2001 May 1;35(5):263-84.
- [8] Hartley RI, Sturm P. Triangulation. *Computer vision and image understanding*. 1997 Nov 1;68(2):146-57.
- [9] Hinton GE, Salakhutdinov RR. Reducing the dimensionality of data with neural networks. *science*. 2006 Jul 28;313(5786):504-7.
- [10] Rumelhart DE, Hinton GE, Williams RJ. Learning representations by back-propagating errors. *nature*. 1986 Oct 9;323(6088):533-6.
- [11] LeCun Y, Bengio Y, Hinton G. Deep learning. *nature*. 2015 May 28;521(7553):436-44.
- [12] LeCun Y, Boser B, Denker JS, et al. Backpropagation applied to handwritten zip code recognition. *Neural computation*. 1989 Dec;1(4):541-51.
- [13] Krizhevsky A, Sutskever I, Hinton GE. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*. 2012;25.
- [14] Szegedy C, Liu W, Jia Y, et al. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition 2015* (pp. 1-9).
- [15] Choy C, Lee J, Ranftl R, et al. High-dimensional convolutional networks for geometric pattern recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition 2020* (pp. 11227-11236).
- [16] Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. In *3rd International Conference on Learning Representations (ICLR 2015) 2015 Apr 10*. Computational and Biological Learning Society.
- [17] He K, Zhang X, Ren S, et al. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition 2016* (pp. 770-778).
- [18] Huang G, Liu Z, Van Der Maaten L, et al. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition 2017* (pp. 4700-4708).
- [19] Ma R, Zhang B, Zhou Y, et al. PID controller-guided attention neural network learning for fast and effective real photographs denoising. *IEEE Transactions on Neural Networks and Learning Systems*. 2021 Jan 15;33(7):3010-23.
- [20] Nair V, Hinton GE. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10) 2010* (pp. 807-814).
- [21] Maas AL, Hannun AY, Ng AY. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml 2013 Jun 16* (Vol. 30, No. 1, p. 3).
- [22] He K, Zhang X, Ren S, et al. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision 2015* (pp. 1026-1034).
- [23] Ioffe S, Szegedy C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning 2015 Jun 1* (pp. 448-456).
- [24] Kingma DP, Ba JL. A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*. 2014 Dec;1412(6).
- [25] Schmidhuber J, Hochreiter S. Long short-term memory. *Neural Comput*. 1997 Nov 15;9(8):1735-80.
- [26] Cho K, Van Merriënboer B, Gulcehre C, et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*. 2014 Jun 3.
- [27] Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need. *Advances in neural information processing systems*. 2017;30.