



Recurrent Vision Transformers for Object Detection with Event Cameras

IT 00 000000 2021111183 000

IT 00 000000 2021114818 000

CONTENTS

01 Before Reviewing..

02 Main Review

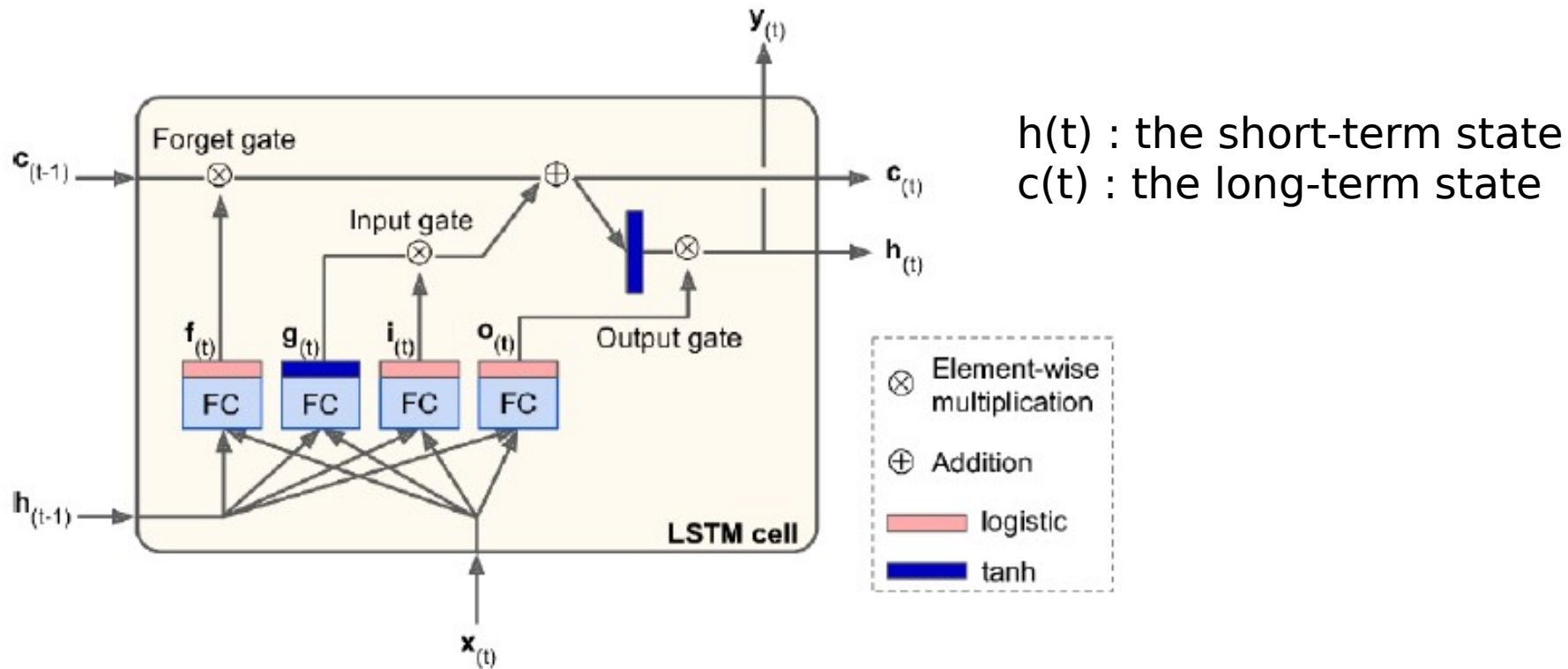
03 Relevance to the subject

01.

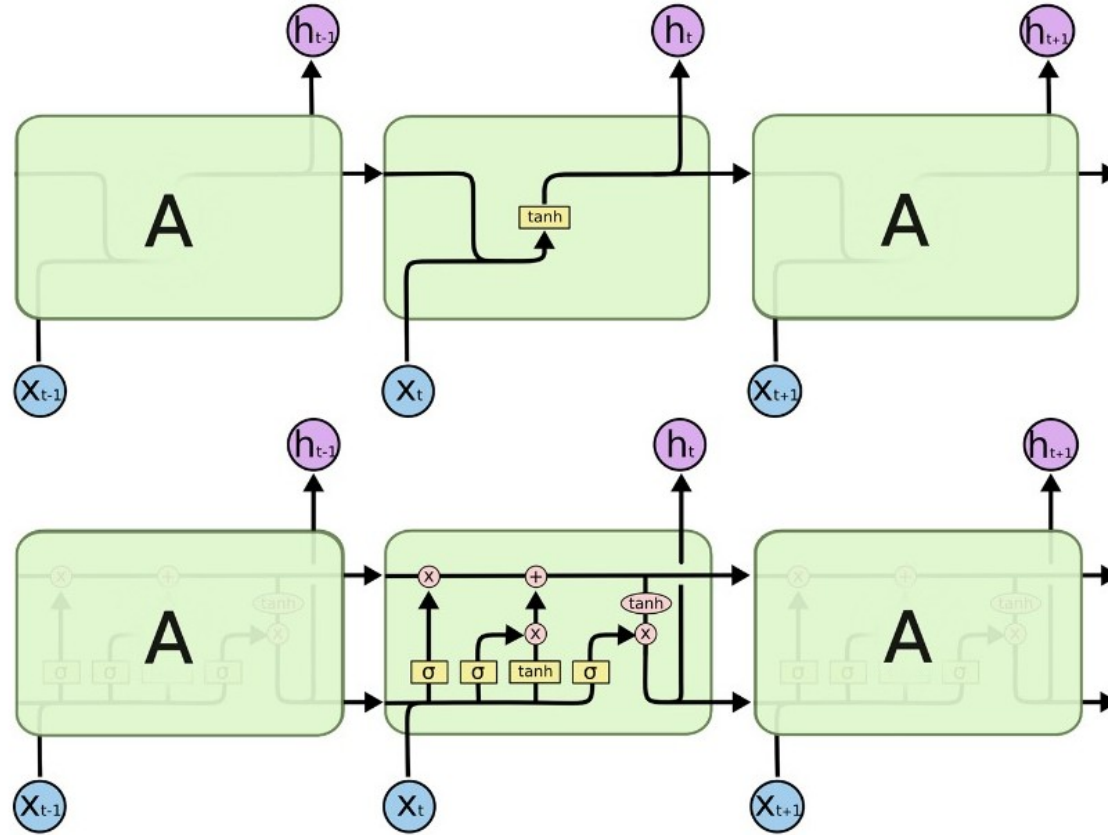
Before Reviewing...

LSTM (Long short-term memory)

RNN $\square\square\square\square$, gate \uparrow $\square\square\square\square\square\square\square\square\square\square$ Vanishing Gradient Problem $\square\square\square\square$.



LSTM (Long short-term memory)



Transformer

“Attention is All you Need”

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

N : sequence length

D : representation dimension

K : kernel size

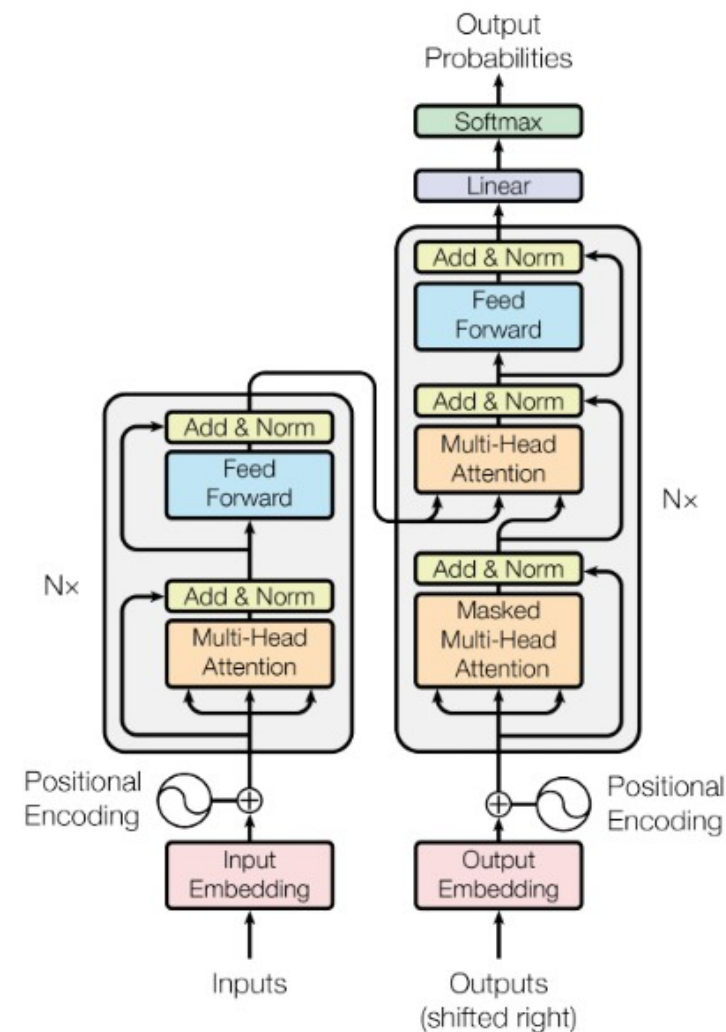
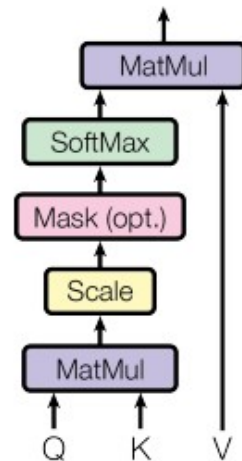


Figure 1: The Transformer - model architecture.

Transformer

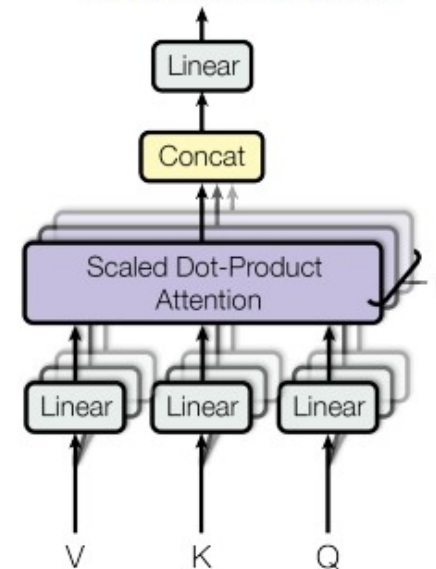
Attention : can be described as mapping a query and a set of key-value pairs to an output

Scaled Dot-Product Attention



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Multi-Head Attention



$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

02.

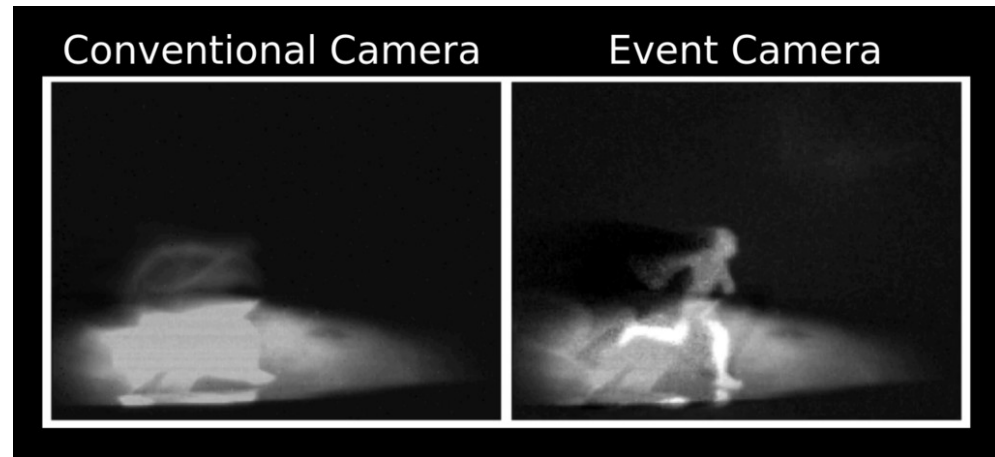
Main Review

1. Introduction

What is **Event Camera**?

□□ □□ □□ □□□ □□□□□□ □□□□ □□□□ □□□□

1. Low latency
2. High dynamic range
3. Strong robustness against motion blur



1. Introduction

Event Cameras vs Traditional Camera(a-frame-based-camera)

	Event Camera	Traditional Camara
absolute intensity information	↓	↑
Change in intensity	↑	↓
Reducing Latency	↑ (submillisecond latency)	↓

Due to latency, traditional camera may come at the cost of missing essential scene details in dynamic scene

1. Introduction

Event Camera란 무엇인가

일반 카메라와 달리, Event Camera는 binary event를 출력한다.
-> 이를 처리하기 위한 detection 및 processing algorithm을 설계해야 한다.

관련 연구 (Related Work)

- GNN(Graph Neural Network)
일반적으로 heavy한 backbone과 ConvLSTM과 같이 expensive한 cell을 사용한다.
-> sparse neural network를 vision backbone로 설계한다.

1. Introduction

□ □□□□ inference time □ performance □ □□□ □□□□ □□ □□□ □□

1. Local and global self-attention
2. Preceding a simple convolution before attention
3. Conv-LSTM => plain LSTM

□□ □□

1. event-based pipelines □ □□□ design □□
2. Simple, composable □ state design
3. State-of-the-art object detection □□ □□□ □□□ □□□ .

2. Related Work

Vision Transformation for Spatio-Temporal Data

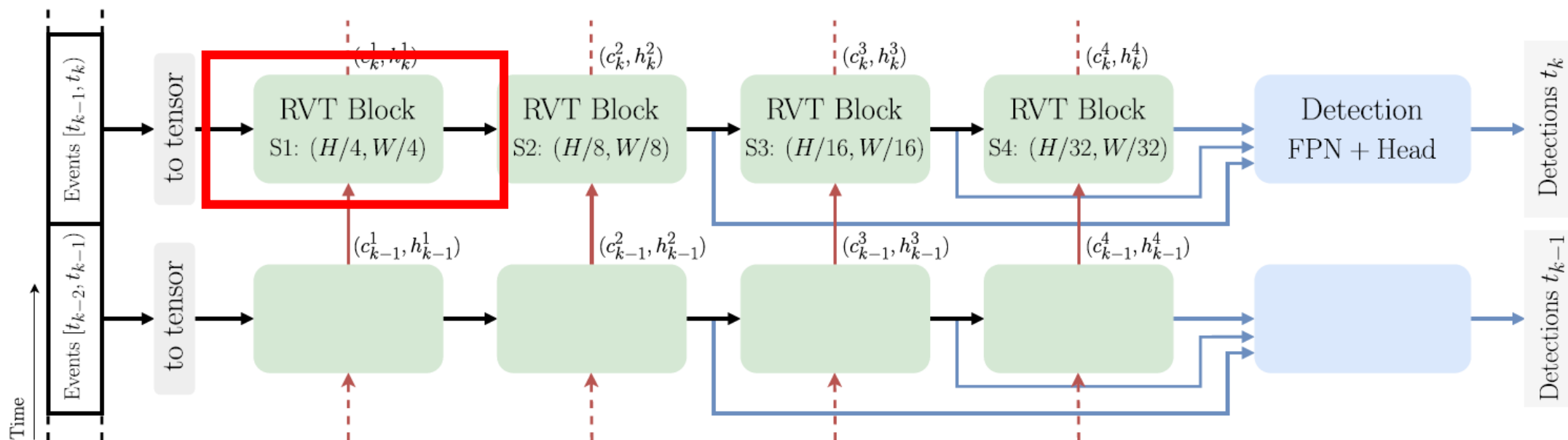
In event-based vision

- classification
- image restructure
- monocular depth estimation

=> Object detection has yet to be investigated

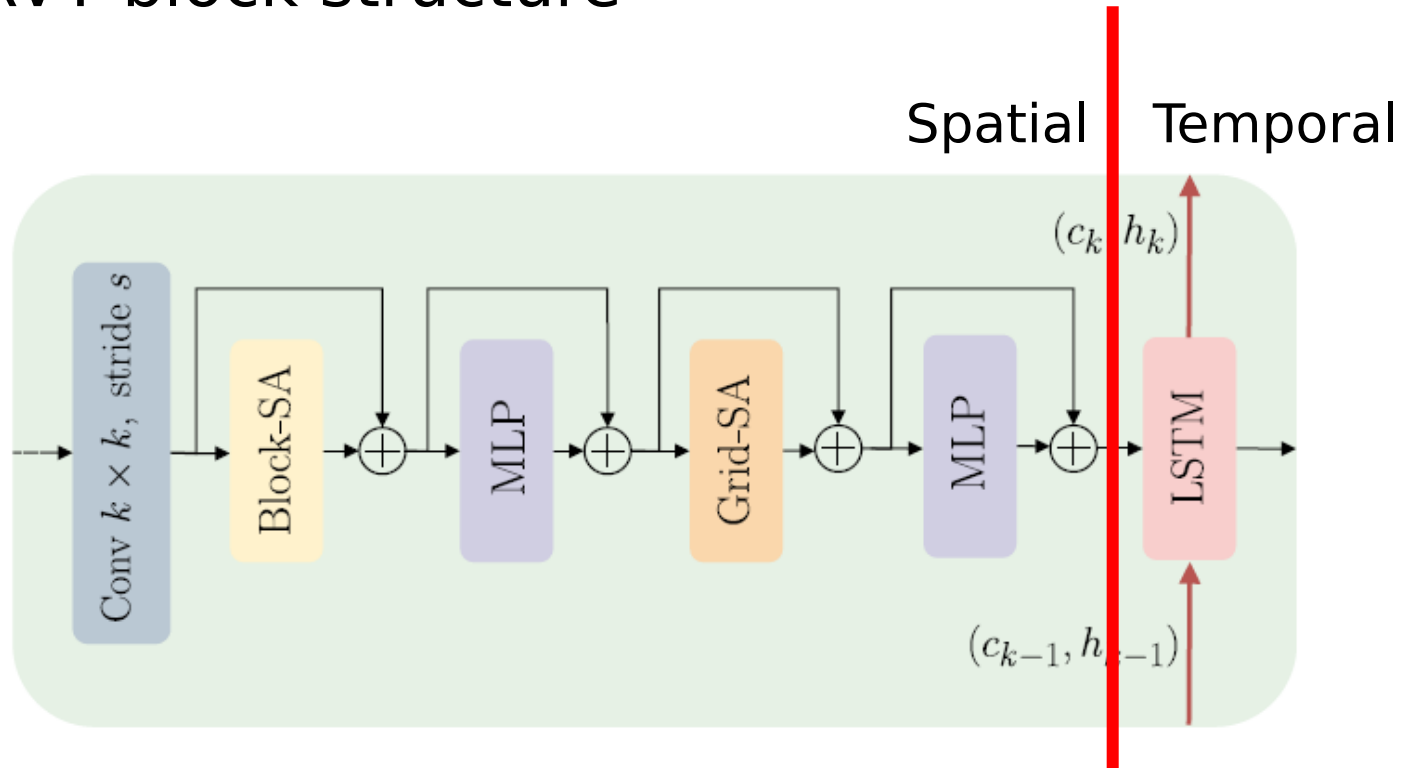
3. Method

Overall Structure



3. Method

RVT block structure



1. Convolution (overlapping)
2. Block Self Attention
3. MLP(Multi Layer Perceptron)
4. Grid Self Attention
5. MLP
6. LSTM

※ Normalization and activation layers are omitted for conciseness

3. Method

Preprocessing step

Input data

$$E(p, \tau, x, y) = \sum_{e_k \in \mathcal{E}} \delta(p - p_k) \delta(x - x_k, y - y_k) \delta(\tau - \tau_k),$$

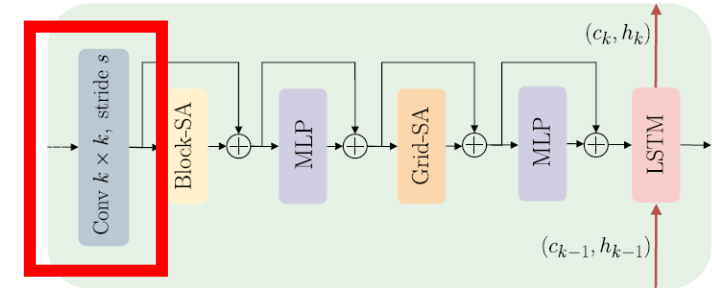
$$\tau_k = \left\lfloor \frac{t_k - t_a}{t_b - t_a} \cdot T \right\rfloor$$

x : width

y : height

p : polarity (□□□)

T : discretization steps of time



(2T, H, W)

2D convolution □ □□□□ □□□ □□

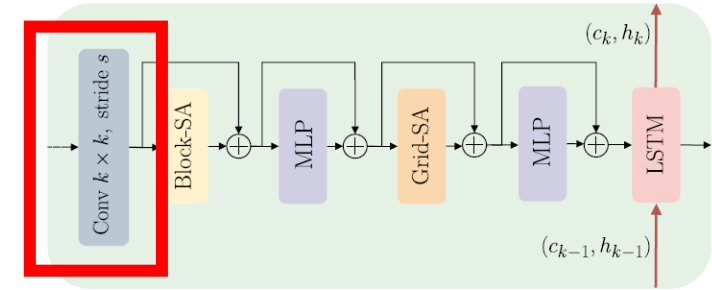
3. Method

Spatial Feature Extraction

1. Convolution (overlapping)

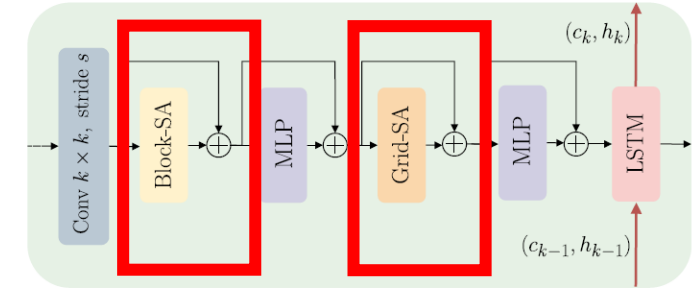
Convolution with overlapping kernels

=> non-overlapping model □□ □□□ □□ □□□ □□□□ .



3. Method

Multi-axis attention self-attention



2. Block Self Attention

Local feature interaction

$$(\frac{H}{P} \times \frac{W}{P}, P \times P, C) P \times P : \text{window size}$$

4. Grid Self Attention

Global feature interaction

$$(G \times G, \frac{H}{G} \times \frac{W}{G}, C)$$

3. Method

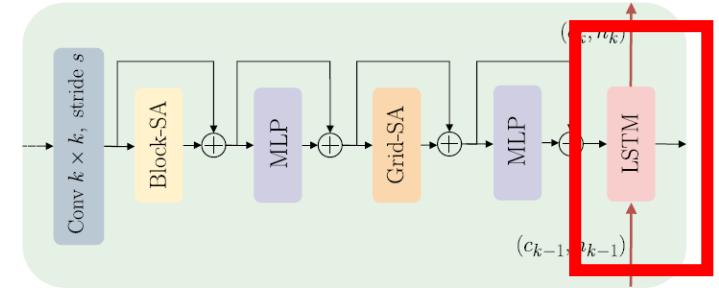
Temporal Feature Extraction

Aggregation with LSTM

Parameter \square $\square\square$ $\square\square\square\square\square$ $\square\square$ Conv-LSTM $\square\square$ Plain LSTM $\square\square$

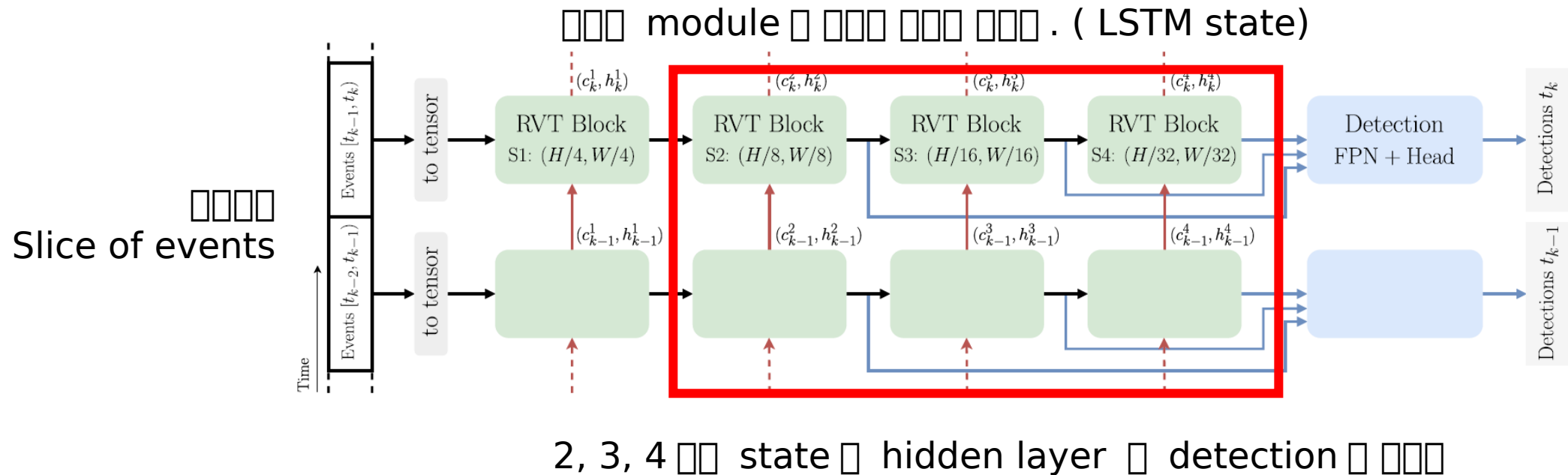
More detail

- \square attention \square MLP $\square\square\square$ $\square\square\square\square\square$ \square , Layer Norm \square $\square\square\square$.
- \square $\square\square\square\square$ Residual connection $\square\square\square$.



3. Method

Hierarchical Multi-Stage Design



4. Experiments

Stage	Size	Kernel	Stride	Channels		
				RVT-T	RVT-S	RVT-B
S1	1/4	7	4	32	48	64
S2	1/8	3	2	64	96	128
S3	1/16	3	2	128	192	256
S4	1/32	3	2	256	384	512

□□ □□ :

- 1.RVT-B (□□ □□): □□ □□ .
- 2.RVT-S (□□ □□): RVT-B □ □□ □□ .
- 3.RVT-T (□□ □□ □□): RVT-B □ □□ □□ □□ .

4. Experiments

4.1 Setup – Implementation Details

1. 실험 환경 :

이 실험은 LayerScale 이 있는 ResNet50 모델을 사용한다 .

2. 실험 방법 :

1. 학습기 ADAM 을 사용하며, 40,000 iteration 까지 학습한다 .
2. OneCycle 을 사용하며, 학습률 스케줄링을 사용한다 .

3. 실험 결과 :

이 실험은 ResNet50 모델을 사용하며, 학습률 스케줄링을 사용하며 Truncated BPTT 를 사용한다 .

4. Experiments

4.1 Setup - Implementation Details

4. **데이터 증강 :**

random horizontal flipping, zooming in and zooming out **랜덤** **랜덤** **랜덤** **랜덤** .

5. **데이터 로딩 :**

50ms **데이터** **데이터** **데이터** , **데이터** 10 **데이터** **데이터** ($T=10$) **데이터** **데이터** **데이터** **데이터** **데이터** .

6. **손실 함수 :**

YOLOX **손실 함수** **손실 함수** , IOU loss, class loss and regression loss 가 **손실 함수** **손실 함수** **손실 함수** **손실 함수** **손실 함수** **손실 함수** .

4. Experiments

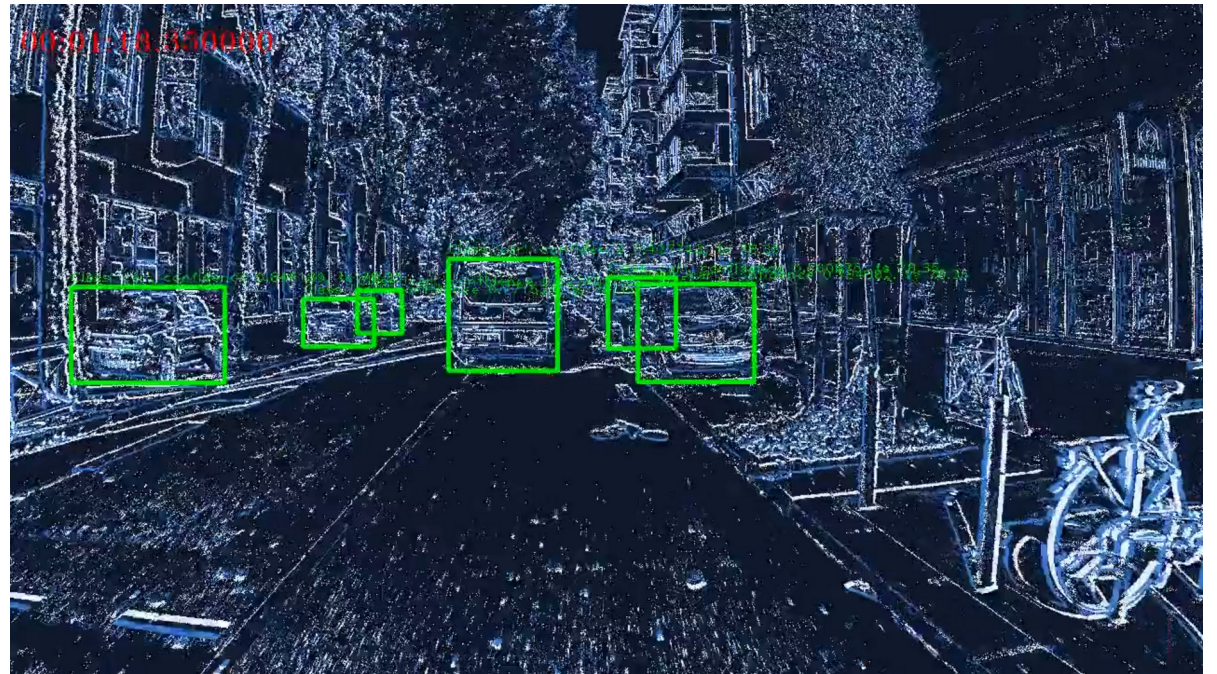
4.1 Setup - Dataset

Gen1 Automotive Detection

□□□□



1 MPx □□□□



4. Experiments

4.2 Ablation Studies – Model Components

Block-type	Gen1		1 Mpx		Params (M)
	mAP	AP ₅₀	mAP	AP ₅₀	
multi-axis	47.6	70.1	46.0	72.3	18.5
Swin	46.7	68.7	44.4	71.7	18.5
ConvNeXt	45.5	65.8	42.3	70.6	18.7

Table 2. **Spatial Aggregation.** Multi-axis attention leads to the best results on both the Gen1 and 1 Mpx dataset.

Conv. kernel type	mAP	AP ₅₀	AP ₇₅	Params (M)
overlapping	47.6	70.1	52.6	18.5
non-overlapping	46.1	68.6	50.5	17.6

Table 3. **Downsampling Strategy.** The usage of overlapping kernels leads to higher performance at the expense of a slight increase in the number of parameters.

4. Experiments

4.2 Ablation Studies – Model Components

LSTM kernel size	mAP	AP ₅₀	AP ₇₅	Params (M)
1×1	47.6	70.1	52.6	18.5
3×3	46.5	69.0	51.4	40.8
3×3 depth-sep	46.3	67.2	51.2	18.6

Table 4. **LSTM kernel size.** Conv-LSTM variants do not outperform the feature specific (1×1) LSTM.

S1	S2	S3	S4	mAP	AP ₅₀	AP ₇₅
				32.0	54.8	31.4
			✓	39.8	63.5	41.6
		✓	✓	44.2	68.4	47.5
	✓	✓	✓	46.9	70.0	50.8
✓	✓	✓	✓	47.6	70.1	52.6

Table 5. **LSTM placement.** LSTM cells contribute to the overall performance even in the early stages.

4. Experiments

4.2 Ablation Studies – Data Augmentations

h-flip	zoom-in	zoom-out	mAP	AP ₅₀	AP ₇₅
			38.1	59.5	41.1
✓			41.6	63.5	45.5
	✓		45.8	67.8	49.8
		✓	44.1	65.7	48.4
<u>✓</u>	<u>✓</u>	<u>✓</u>	47.6	70.1	52.6

Table 7. **Data Augmentation.** Data augmentation consistently improves the results.

4. Experiments

4.3 Benchmark Comparisons

Method	Backbone	Detection Head	Gen1		1 Mpx		Params (M)
			mAP	Time (ms)	mAP	Time (ms)	
NVS-S [27]	GNN	YOLOv1 [40]	8.6	-	-	-	0.9
Asynet [34]	Sparse CNN	YOLOv1	14.5	-	-	-	11.4
AEGNN [43]	GNN	YOLOv1	16.3	-	-	-	20.0
Spiking DenseNet [10]	SNN	SSD [30]	18.9	-	-	-	8.2
Inception + SSD [19]	CNN	SSD	30.1	19.4	34.0	45.2	> 60*
RRC-Events [7]	CNN	YOLOv3 [41]	30.7	21.5	34.3	46.4	> 100*
MatrixLSTM [6]	RNN + CNN	YOLOv3	31.0	-	-	-	61.5
YOLOv3 Events [20]	CNN	YOLOv3	31.2	22.3	34.6	49.4	> 60*
RED [38]	CNN + RNN	SSD	40.0	16.7	43.0	39.3	24.1
ASTMNet [26]	(T)CNN + RNN	SSD	46.7	35.6	48.3	72.3	> 100*
RVT-B (ours)	Transformer + RNN	YOLOX [15]	47.2	10.2 (3.7)	<u>47.4</u>	11.9 (6.1)	18.5
RVT-S (ours)	Transformer + RNN	YOLOX	46.5	9.5 (3.0)	44.1	10.1 (5.0)	9.9
RVT-T (ours)	Transformer + RNN	YOLOX	44.1	9.4 (2.3)	41.5	9.5 (3.5)	4.4

4. Experiments

4.3 Benchmark Comparisons

- **Gen1** :

- Gen1 **47.2 mAP**, 1 MPx **47.4 mAP** **5** **0.1** **0.1** **0.1**

- **ASTMNet**:

- ASTMNet backbone 가 **4** **0.1** **0.1** **0.1**

- **RED** :

- RED **4.1** **0.1** **0.1** **0.1** Gen1 **4.1** **0.1** **0.1** **0.1** **mAP** 가 **7.2** **0.1** ,
- 1 MPx **4.4** **0.1** **0.1**

- **Tiny** :

- Gen1 **4.1** **0.1** **0.1** **0.1** RED **4.1** **0.1** **0.1** **0.1** **mAP** **5** **0.1** **0.1** **0.1**

4. Experiments

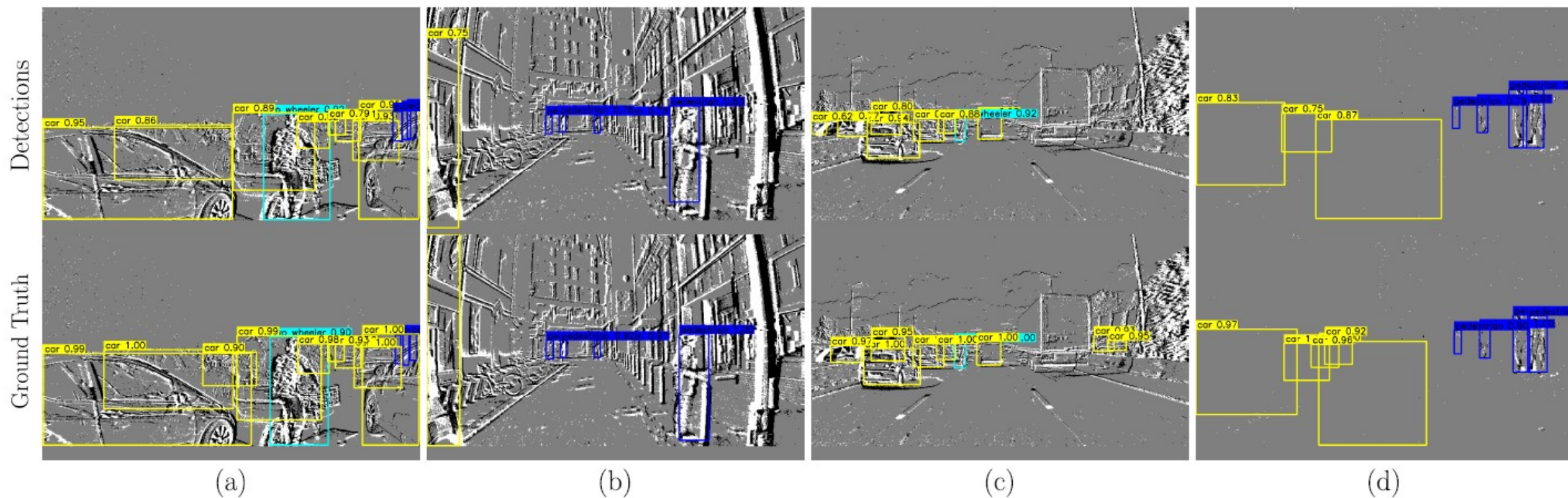


Figure 4. **Predictions on the 1 Mpx dataset.** All examples are thematically picked to illustrate the behaviour of the model in different scenarios. (d) shows a scenario in which the model can still partially detect objects in absence of events due to the temporal memory.

6. Conclusion

요약 :

- 기존 방법론에서 사용된 Backbone 구조를 변경 .
- 기존 방법론에서 사용된 Attention 구조를 변경 .

주요 특징 :

Convolution prior
local- and sparse global attention
recurrent feature aggregation

결과 :

- RVT 구조를 사용하여 기존 방법론에서 사용된 Attention 구조를 변경 .

주요 성과 :

- 기존 방법론에서 사용된 Attention 구조를 변경 , 기존 방법론에서 사용된 low -latency 구조를 변경 .

03.

Relevance to the subject

Relevance to the subject

Optimizer : Adam

Batch Strategy: BPTT , Truncated BPTT

Precision: MAP

Model LSTM, Transformer, CNN

Etc..

Overfitting □ □□□ □□ □□□ □□□□ □□

MLP □□

Residual □□ □□



https://m.hanbit.co.kr/channel/category/category_view.html?cms_code=CMS6074576268
<https://wikidocs.net/31379>
https://en.wikipedia.org/wiki/Event_camera
https://ko.wikipedia.org/wiki/%EB%AA%A8%EC%85%98_%EB%B8%94%EB%9F%AC

