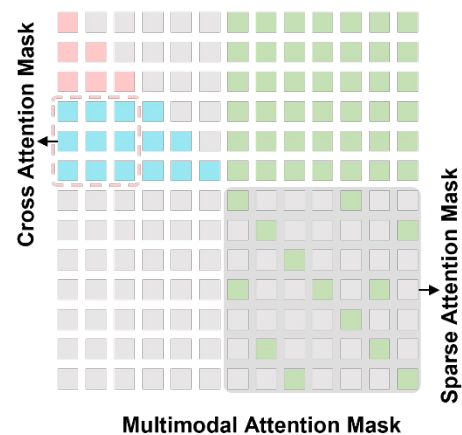
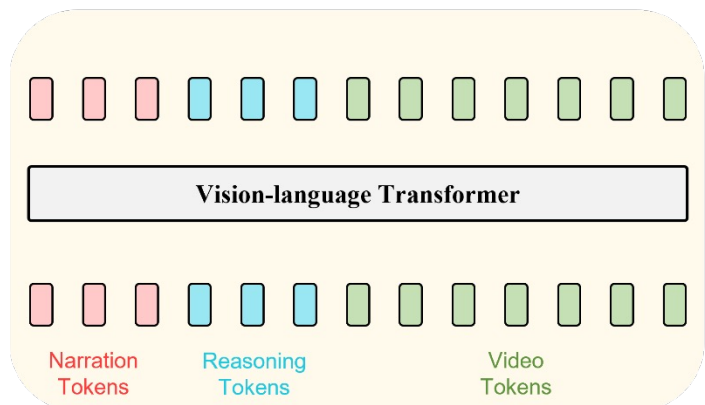
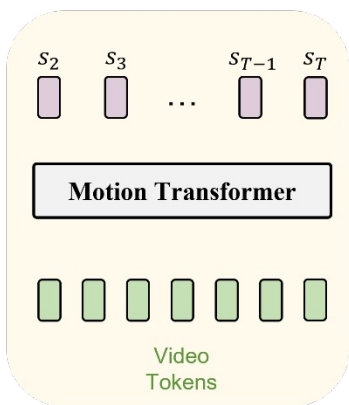
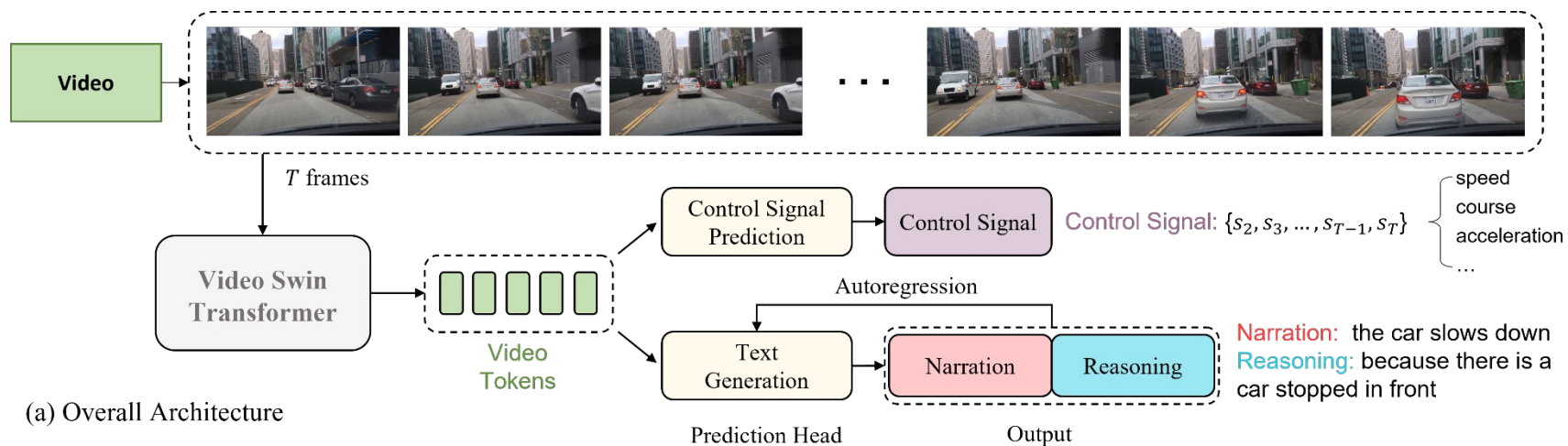


# Implementation



# Result

Method	Narration			Reasoning		
	B4	C	M	B4	C	M
S2VT [42]	30.2	179.8	27.5	6.3	53.4	11.2
S2VT++ [42]	27.1	157.0	26.4	5.8	52.7	10.9
SAA [56]	31.8	214.8	29.1	7.1	66.1	12.2
WAA [56]	32.3	215.8	29.2	7.3	69.5	12.2
Ours	<b>34.6</b>	<b>247.5</b>	<b>30.6</b>	<b>11.4</b>	<b>102.6</b>	<b>15.2</b>

Method	Narration	Reasoning	Full sentence
SAA [56]	90.8%	62.4%	-
WAA [56]	<b>93.5%</b>	66.0%	-
Ours	90.0%	<b>90.3%</b>	82.7%

Comparison with sota results: machine evaluation (top) and manual evaluation (bottom)

# Ablation Experiment

Method	Narration				Reasoning			
	B4	C	M	R	B4	C	M	R
Single	33.2	238.9	29.7	62.0	8.6	89.7	14.1	31.4
Single+	33.9	<b>248.3</b>	30.5	<b>63.1</b>	9.3	97.2	14.6	31.5
Ours	<b>34.6</b>	247.5	<b>30.6</b>	62.8	<b>11.4</b>	<b>102.6</b>	<b>15.2</b>	<b>32.0</b>

Comparison of single-task and multi-task experiments

# Ablation Experiment

Signals		Narration			Reasoning		
Speed	Course	C	M	R	C	M	R
✓		232.0	29.9	61.5	88.0	15.1	31.0
	✓	218.2	29.3	61.2	88.6	14.1	30.6
✓	✓	<b>247.5</b>	<b>30.6</b>	<b>62.8</b>	<b>102.6</b>	<b>15.2</b>	<b>32.0</b>

Effect of different control signals on caption results

# Ablation Experiment

Method	Narration				Reasoning				Cost(min)
	B4	C	M	R	B4	C	M	R	
2	33.4	227.7	28.7	61.0	8.7	62.9	15.1	29.8	294
4	32.9	225.7	29.0	60.9	9.9	81.3	14.9	31.1	382
8	32.6	236.1	29.3	61.8	8.4	83.7	13.4	30.6	447
16	32.5	231.0	29.5	61.9	8.7	91.5	13.8	32.0	528
32	<b>34.6</b>	<b>247.5</b>	<b>30.6</b>	<b>62.8</b>	<b>11.4</b>	<b>102.6</b>	<b>15.2</b>	<b>32.0</b>	797

Effect of video frame number on the result

# Ablation Experiment

$$c_{\sigma} = \begin{cases} 1, & -\sigma < \hat{c} - c < \sigma \\ 0, & \textit{otherwise} \end{cases}$$

Method	Course						Speed					
	RMSE(degree)↓	$A_{0.1} \uparrow$	$A_{0.5} \uparrow$	$A_{1.0} \uparrow$	$A_{5.0} \uparrow$	$A_{10.0} \uparrow$	RMSE(m/s)↓	$A_{0.1} \uparrow$	$A_{0.5} \uparrow$	$A_{1.0} \uparrow$	$A_{5.0} \uparrow$	$A_{10.0} \uparrow$
Single	<b>6.3</b>	8.3	84.7	<b>90.5</b>	97.2	98.7	3.4	5.0	25.5	37.8	86.8	98.7
Ours	6.4	<b>62.2</b>	<b>85.5</b>	89.9	<b>97.2</b>	<b>98.8</b>	<b>2.5</b>	<b>11.1</b>	<b>28.1</b>	<b>45.3</b>	<b>94.3</b>	<b>99.5</b>

Effect of Multitasking on Control Signal Prediction

# Visualization

