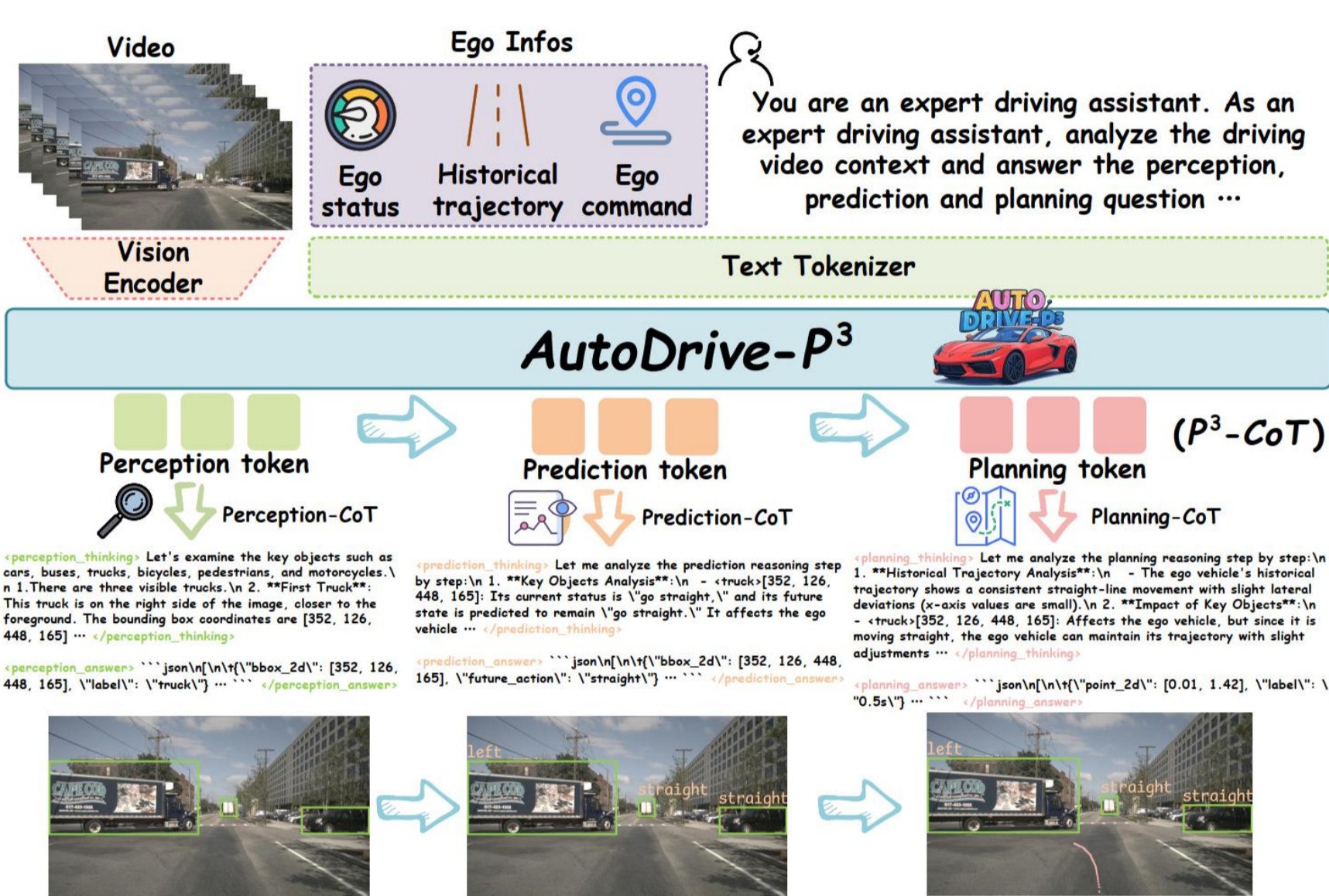


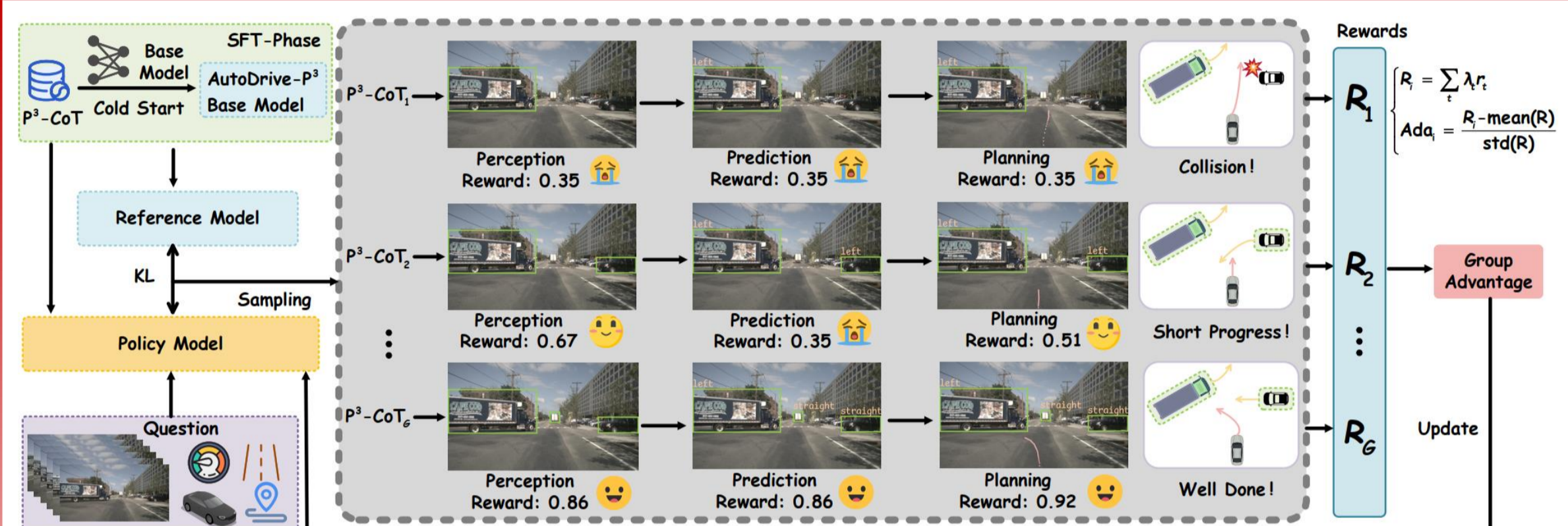
Yuqi Ye^{†,1}, Zijian Zhang^{†,1}, Junhong Lin¹, Shangkun Sun¹, Changhao Peng¹, Wei Gao^{*,1}¹School of Electronic and Computer Engineering, Peking University[†]Core Authors, * ✉ gaoweit262@pku.edu.cn北京大学信息工程学院
School of Electronic and Computer Engineering
Peking University

Abstract

Vision-language models (VLMs) are increasingly being adopted for end-to-end autonomous driving systems due to their exceptional performance in handling long-tail scenarios. However, current VLM-based approaches suffer from two major limitations: 1) Some VLMs directly output planning results without chain-of-thought (CoT) reasoning, bypassing crucial perception and prediction stages which creates a significant domain gap and compromises decision making capability; 2) Other VLMs can generate outputs for perception, prediction, and planning tasks but employ a fragmented decision-making approach where these modules operate separately, leading to a significant lack of synergy that undermines true planning performance. To address these limitations, we propose AutoDrive-P³, a novel framework that seamlessly integrates Perception, Prediction, and Planning through structured reasoning. We introduce the P³-CoT dataset to facilitate coherent reasoning and propose P³-GRPO, a hierarchical reinforcement learning algorithm that provides progressive supervision across all three tasks. Specifically, AutoDrive-P³ progressively generates CoT reasoning and answers for perception, prediction, and planning, where perception provides essential information for subsequent prediction and planning, while both perception and prediction collectively contribute to the final planning decisions, enabling safer and more interpretable autonomous driving. Additionally, to balance inference efficiency with performance, we introduce dual thinking modes: detailed thinking and fast thinking. Extensive experiments on both open-loop (nuScenes) and closed-loop (NAVSIMv1/v2) benchmarks demonstrate that our approach achieves state-of-the-art performance in planning tasks.

AutoDrive-P³ Pipeline

AutoDrive-P³ processes video and ego vehicle data through structured Perception-Prediction-Planning Chain-of-Thought (P³-CoT) reasoning, generating interpretable step-by-step rationale and structured outputs for perception, prediction, and planning.

P³-GRPO Pipeline

We first cold start the base model using P³-CoT to make up for the gap between VLM and autonomous driving and learn the CoT answer format. Next we further enhance the VLM's reasoning capability across all three stages by applying the GRPO algorithm to the perception, prediction, and planning modules collectively, yielding the P³-GRPO algorithm. Our approach employs a multicomponent reward function to guide the policy model toward generating accurate, coherent, and well-structured outputs through coordinated reinforcement learning across these cognitive layers.

Results

Table 1: Performance comparison on nuScenes Benchmark.

Method	L2 (m) ↓				Collision (%) ↓				VLM
	1s	2s	3s	Avg.	1s	2s	3s	Avg.	
Non-Autoregressive Methods									
ST-P3 (Hu et al. [2022])	1.33	2.11	2.90	2.11	0.23	0.62	1.27	0.71	-
VAD (Jiang et al. [2023])	0.17	0.34	0.60	0.37	0.07	0.10	0.24	0.14	-
Ego-MLP (Li et al. [2024d])	0.46	0.76	1.12	0.78	0.21	0.35	0.58	0.38	-
UniAD (Hu et al. [2023])	0.44	0.67	0.96	0.69	0.04	0.08	0.23	0.12	-
InsightDrive (Song et al. [2025])	0.23	0.41	0.68	0.44	0.09	0.10	0.27	0.15	-
Autoregressive Methods									
GPT-Driver (Mao et al. [2023])	0.20	0.40	0.70	0.44	0.04	0.12	0.36	0.17	GPT-3.5
DriveVLM (Tian et al. [2024])	0.18	0.34	0.68	0.40	0.10	0.22	0.45	0.27	Qwen2-VL-7B
OpenEMMA (Xing et al. [2025])	1.45	3.21	3.76	2.81	-	-	-	-	Qwen2-VL-7B
RDA-Driver (Huang et al. [2024a])	0.17	0.37	0.69	0.40	0.01	0.05	0.26	0.10	LLaVa-7B
OmniDrive (Wang et al. [2024])	0.14	0.29	0.55	0.33	0.01	0.04	0.27	0.11	LLaVa-7B
OpenDriveVLA (Zhou et al. [2025a])	0.14	0.30	0.55	0.33	0.02	0.07	0.22	0.10	Qwen2.5-VL-3B
AutoVLA (Zhou et al. [2025b])	0.25	0.46	0.73	0.48	0.07	0.07	0.26	0.13	Qwen2.5-VL-3B
AutoDrive-R ² (Yuan et al. [2025])	0.35	0.49	0.62	0.49	-	-	-	-	Qwen2.5-VL-3B
AutoDrive-P ³ (Ours-Detailed)	0.15	0.30	0.54	0.33	0.00	0.02	0.15	0.06	Qwen2.5-VL-3B
AutoDrive-P ³ (Ours-Fast)	0.16	0.31	0.56	0.34	0.00	0.04	0.20	0.08	Qwen2.5-VL-3B

Table 2: Performance comparison on NAVSIMv1 benchmark.

Method	Image	Lidar	NC↑	DAC↑	EP↑	TTC↑	Comf↑	PDMS↑
Human	X	X	100.0	100.0	87.5	100.0	99.9	94.8
Constant Velocity	X	X	69.9	58.8	49.3	49.3	100.0	21.6
Ego Status MLP	X	X	93.0	77.3	62.8	83.6	100.0	65.6
VADv2 (Weng et al. [2024])	✓	X	97.9	91.7	77.6	92.9	100.0	83.0
UniAD (Hu et al. [2023])	✓	X	97.8	91.9	78.8	92.9	100.0	83.4
LTF (Prakash et al. [2021])	✓	X	97.4	92.8	79.0	92.4	100.0	83.8
TransFuser (Prakash et al. [2021])	✓	✓	97.7	92.8	79.2	92.8	100.0	84.0
PARA-Drive (Weng et al. [2024])	✓	X	97.9	92.4	79.3	93.0	99.8	84.0
LAW (Li et al. [2024a])	✓	✓	96.4	95.4	81.7	88.7	99.9	84.6
DRAMA (Yuan et al. [2024])	✓	✓	98.0	93.1	80.1	94.8	100.0	85.5
Hydra-MDP (Li et al. [2024b])	✓	✓	98.3	96.0	78.7	94.6	100.0	86.5
DiffusionDrive (Liao et al. [2025])	✓	✓	98.2	96.2	82.2	94.7	100.0	88.1
WoTE (Li et al. [2025])	✓	✓	98.5	96.8	81.9	94.9	99.9	88.3
AutoDrive-P ³ (Ours-Detailed)	✓	X	99.1	97.4	84.8	96.5	100.0	90.6
AutoDrive-P ³ (Ours-Fast)	✓	X	98.9	97.7	83.7	96.6	99.9	90.2

Table 5: Ablation study on different training setting on nuScenes benchmark.

Method	Group Size	History Traj.	Sensor Type	L2 (m) ↓				Collision (%) ↓			
				1s	2s	3s	Avg.	1s	2s	3s	Avg.
Ablation 1	4	✓	Video	0.17	0.32	0.65	0.38	0.01	0.06	0.30	0.13
Ablation 2	8	X	Video	0.17	0.33	0.68	0.39	0.02	0.07	0.33	0.14
Ablation 3	8	✓	Image	0.16	0.32	0.61	0.36	0.01	0.05	0.26	0.12
P ³ -GRPO	8	✓	Video	0.15	0.30	0.54	0.33	0.00	0.02	0.15	0.06

Table 3: Performance comparison on NAVSIMv2 benchmark.

Method	NC↑	DAC↑	DDC↑	TLC↑	EP↑	TTC↑	LK↑	HC↑	EC↑	EPDMS↑ False / True
Human	100.0	100.0	99.8	100.0	87.4	100.0	100.0	98.1	90.1	90.3 / 94.5
Ego Status MLP	93.1	77.9	92.7	99.6	86.0	91.5	89.4	98.3	85.4	64.0 / -
Transfuser (Prakash et al. [2021])	96.9	89.9	97.8	99.7	87.1	95.4	92.7	98.3	87.2	76.7 / 84.0
HydraMDP++ (Li et al. [2024b])	97.2	97.5	99.4	99.6	83.1	96.5	94.4	98.2	70.9	81.4 / -
DiffusionDrive (Liao et al. [2025])	98.2	96.2	99.5	99.8	87.4	97.3	96.9	98.4	87.7	84.7 / 88.2
WoTE (Li et al. [2025])	98.5	96.8	98.8	99.8	86.1	97.9	95.5	98.3	82.9	84.2 / 87.7
AutoDrive-P ³ (Ours-Detailed)	99.1	97.4	99.2	99.8	88.0	98.7	96.3	98.3	85.5	86.2 / 89.9
AutoDrive-P ³ (Ours-Fast)	98.9	97.6	98.9	99.8	86.8	98.5	95.4	98.3	80.6	85.2 / 88.7

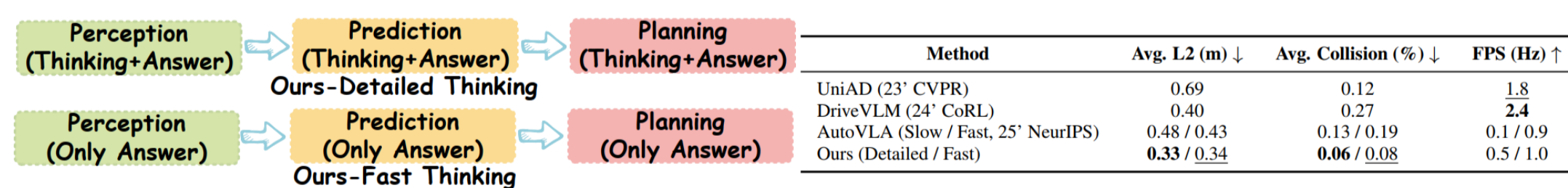


Figure 5: Dual thinking modes and running time on nuScenes Benchmark.

```

OURS-FAST:
<think_perception></think_perception>
<answer_perception> json\n{\n\t\bbox_2d\n": [], \n\tlabel\n": []\n}\n</answer_perception>
<think_prediction></think_prediction>
<answer_prediction> json\n{\n\t\bbox_2d\n": [], \n\tfuture_position\n": []\n}\n</answer_prediction>
<think_planning></think_planning>
<answer_planning> json\n{\n\t\tx_y_radian\n": [1.35, 0.0, 0.0], \n\tlabel\n": "\0.5s\n", \n\t\tx_y_radian\n": [3.15, -0.02, -0.02], \n\tlabel\n": "\1.0s\n", \n\t\tx_y_radian\n": [5.25, -0.11, -0.08], \n\tlabel\n": "\1.5s\n", \n\t\tx_y_radian\n": [7.55, -0.41, -0.19], \n\tlabel\n": "\2.0s\n", \n\t\tx_y_radian\n": [9.89, -1.02, -0.35], \n\tlabel\n": "\2.5s\n", \n\t\tx_y_radian\n": [12.1, -2.05, -0.55], \n\tlabel\n": "\3.0s\n", \n\t\tx_y_radian\n": [13.98, -3.51, -0.78], \n\tlabel\n": "\3.5s\n", \n\t\tx_y_radian\n": [15.42, -5.35, -1.02], \n\tlabel\n": "\4.0s\n", \n\t\tx_y_radian\n": [17.35, -7.35, -1.42], \n\tlabel\n": "\4.5s\n", \n\t\tx_y_radian\n": [19.35, -9.6, -1.95], \n\tlabel\n": "\5.0s\n", \n\t\tx_y_radian\n": [21.45, -12.15, -2.65], \n\tlabel\n": "\5.5s\n", \n\t\tx_y_radian\n": [23.65, -15.05, -3.55], \n\tlabel\n": "\6.0s\n", \n\t\tx_y_radian\n": [25.95, -18.3, -4.65], \n\tlabel\n": "\6.5s\n", \n\t\tx_y_radian\n": [28.35, -21.85, -5.95], \n\tlabel\n": "\7.0s\n", \n\t\tx_y_radian\n": [30.85, -26.7, -7.45], \n\tlabel\n": "\7.5s\n", \n\t\tx_y_radian\n": [33.45, -31.85, -9.15], \n\tlabel\n": "\8.0s\n", \n\t\tx_y_radian\n": [36.15, -38.2, -11.05], \n\tlabel\n": "\8.5s\n", \n\t\tx_y_radian\n": [38.95, -45.85, -13.15], \n\tlabel\n": "\9.0s\n", \n\t\tx_y_radian\n": [41.85, -54.75, -15.45], \n\tlabel\n": "\9.5s\n", \n\t\tx_y_radian\n": [44.85, -64.95, -17.95], \n\tlabel\n": "\10.0s\n", \n\t\tx_y_radian\n": [47.95, -76.45, -20.65], \n\tlabel\n": "\10.5s\n", \n\t\tx_y_radian\n": [51.15, -89.25, -23.55], \n\tlabel\n": "\11.0s\n", \n\t\tx_y_radian\n": [54.45, -103.35, -26.65], \n\tlabel\n": "\11.5s\n", \n\t\tx_y_radian\n": [57.85, -118.75, -29.95], \n\tlabel\n": "\12.0s\n", \n\t\tx_y_radian\n": [61.35, -135.45, -33.45], \n\tlabel\n": "\12.5s\n", \n\t\tx_y_radian\n": [64.95, -153.45, -37.15], \n\tlabel\n": "\13.0s\n", \n\t\tx_y_radian\n": [68.65, -172.75, -41.05], \n\tlabel\n": "\13.5s\n", \n\t\tx_y_radian\n": [72.45, -193.35, -45.15], \n\tlabel\n": "\14.0s\n", \n\t\tx_y_radian\n": [76.35, -215.25, -49.45], \n\tlabel\n": "\14.5s\n", \n\t\tx_y_radian\n": [80.35, -238.45, -53.95], \n\tlabel\n": "\15.0s\n", \n\t\tx_y_radian\n": [84.45, -262.85, -58.65], \n\tlabel\n": "\15.5s\n", \n\t\tx_y_radian\n": [88.65, -288.45, -63.55], \n\tlabel\n": "\16.0s\n", \n\t\tx_y_radian\n": [92.95, -315.25, -68.65], \n\tlabel\n": "\16.5s\n", \n\t\tx_y_radian\n": [97.35, -343.35, -73.95], \n\tlabel\n": "\17.0s\n", \n\t\tx_y_radian\n": [101.85, -372.65, -79.45], \n\tlabel\n": "\17.5s\n", \n\t\tx_y_radian\n": [106.45, -403.15, -85.15], \n\tlabel\n": "\18.0s\n", \n\t\tx_y_radian\n": [111.15, -434.85, -91.05], \n\tlabel\n": "\18.5s\n", \n\t\tx_y_radian\n": [115.95, -467.75, -97.15], \n\tlabel\n": "\19.0s\n", \n\t\tx_y_radian\n": [120.85, -501.85, -103.45], \n\tlabel\n": "\19.5s\n", \n\t\tx_y_radian\n": [125.85, -537.15, -109.95], \n\tlabel\n": "\20.0s\n", \n\t\tx_y_radian\n": [130.95, -573.65, -116.65], \n\tlabel\n": "\20.5s\n", \n\t\tx_y_radian\n": [136.15, -611.35, -123.55], \n\tlabel\n": "\21.0s\n", \n\t\tx_y_radian\n": [141.45, -650.25, -130.65], \n\tlabel\n": "\21.5s\n", \n\t\tx_y_radian\n": [146.85, -690.35, -137.95], \n\tlabel\n": "\22.0s\n", \n\t\tx_y_radian\n": [152.35, -731.65, -145.45], \n\tlabel\n": "\22.5s\n", \n\t\tx_y_radian\n": [157.95, -774.15, -153.15], \n\tlabel\n": "\23.0s\n", \n\t\tx_y_radian\n": [163.65, -817.85, -161.05], \n\tlabel\n": "\23.5s\n", \n\t\tx_y_radian\n": [169.45, -862.75, -169.15], \n\tlabel\n": "\24.0s\n", \n\t\tx_y_radian\n": [175.35, -908.85, -177.45], \n\tlabel\n": "\24.5s\n", \n\t\tx_y_radian\n": [181.35, -956.15, -185.95], \n\tlabel\n": "\25.0s\n", \n\t\tx_y_radian\n": [187.45, -1004.65, -194.65], \n\tlabel\n": "\25.5s\n", \n\t\tx_y_radian\n": [193.65, -1054.35, -203.55], \n\tlabel\n": "\26.0s\n", \n\t\tx_y_radian\n": [199.95, -1105.25, -212.65], \n\tlabel\n": "\26.5s\n", \n\t\tx_y_radian\n": [206.35, -1157.35, -221.95], \n\tlabel\n": "\27.0s\n", \n\t\tx_y_radian\n": [212.85, -1210.65, -231.45], \n\tlabel\n": "\27.5s\n", \n\t\tx_y_radian\n": [219.45, -1265.15, -241.15], \n\tlabel\n": "\28.0s\n", \n\t\tx_y_radian\n": [226.15, -1320.85, -251.05], \n\tlabel\n": "\28.5s\n", \n\t\tx_y_radian\n": [232.95, -1377.75, -261.15], \n\tlabel\n": "\29.0s\n", \n\t\tx_y_radian\n": [239.85, -1435.85, -271.45], \n\tlabel\n": "\29.5s\n", \n\t\tx_y_radian\n": [246.85, -1495.15, -281.95], \n\tlabel\n": "\30.0s\n", \n\t\tx_y_radian\n": [253.95, -1555.65, -292.65], \n\tlabel\n": "\30.5s\n", \n\t\tx_y_radian\n": [261.15, -1617.35, -303.55], \n\tlabel\n": "\31.0s\n", \n\t\tx_y_radian\n": [268.45, -1680.25, -314.65], \n\tlabel\n": "\31.5s\n", \n\t\tx_y_radian\n": [275.85, -1744.35, -325.95], \n\tlabel\n": "\32.0s\n", \n\t\tx_y_radian\n": [283.35, -1809.65, -337.45], \n\tlabel\n": "\32.5s\n", \n\t\tx_y_radian\n": [290.95, -1876.15, -349.15], \n\tlabel\n": "\33.0s\n", \n\t\tx_y_radian\n": [298.65, -1943.85, -361.05], \n\tlabel\n": "\33.5s\n", \n\t\tx_y_radian\n": [306.45, -2012.75, -373.15], \n\tlabel\n": "\34.0s\n", \n\t\tx_y_radian\n": [314.35, -2082.85, -385.45], \n\tlabel\n": "\34.5s\n", \n\t\tx_y_radian\n": [322.35, -2154.15, -397.95], \n\tlabel\n": "\35.0s\n", \n\t\tx_y_radian\n": [330.45, -2226.65, -410.65], \n\tlabel\n": "\35.5s\n", \n\t\tx_y_radian\n": [338.65, -2300.35, -423.55], \n\tlabel\n": "\36.0s\n", \n\t\tx_y_radian\n": [346.95, -2375.25, -436.65], \n\tlabel\n": "\36.5s\n", \n\t\tx_y_radian\n": [355.35, -2451.35, -449.95], \n\tlabel\n": "\37.0s\n", \n\t\tx_y_radian\n": [363.85, -2528.65, -463.45], \n\tlabel\n": "\37.5s\n", \n\t\tx_y_radian\n": [372.45, -2607.15, -477.15], \n\tlabel\n": "\38.0s\n", \n\t\tx_y_radian\n": [381.15, -2686.85, -491.05], \n\tlabel\n": "\38.5s\n", \n\t\tx_y_radian\n": [390.05, -2767.75, -505.15], \n\tlabel\n": "\39.0s\n", \n\t\tx_y_radian\n": [399.15, -2849.85, -519.45], \n\tlabel\n": "\39.5s\n", \n\t\tx_y_radian\n": [408.35, -2933.15, -533.95], \n\tlabel\n": "\40.0s\n", \n\t\tx_y_radian\n": [417.65, -3017.65, -548.65], \n\tlabel\n": "\40.5s\n", \n\t\tx_y_radian\n": [427.15, -3103.35, -563.55], \n\tlabel\n": "\41.0s\n", \n\t\tx_y_radian\n": [436.85, -3190.25, -578.65], \n\tlabel\n": "\41.5s\n", \n\t\tx_y_radian\n": [446.65, -3278.35, -593.95], \n\tlabel\n": "\42.0s\n", \n\t\tx_y_radian\n": [456.65, -3367.65, -609.45], \n\tlabel\n": "\42.5s\n", \n\t\tx_y_radian\n": [466.85, -3458.15, -625.15], \n\tlabel
```