

Policy-Guided World Model Planning for Language-Conditioned Visual Navigation

Anonymous CVPR submission

Paper ID ****

Abstract

001 Navigating to a visually specified goal given natural lan- 036
002 guage instructions remains a fundamental challenge in em- 037
003 bodied AI. Existing approaches either rely on reactive poli- 038
004 cies that struggle with long-horizon planning, or employ 039
005 world models that suffer from poor action initialization 040
006 in high-dimensional spaces. We present PiJEPa, a two- 041
007 stage framework that combines the strengths of learned 042
008 navigation policies with latent world model planning for 043
009 instruction-conditioned visual navigation. In the first stage, 044
010 we finetune an Octo-based generalist policy, augmented 045
011 with a frozen pretrained vision encoder (DINOv2 or V- 046
012 JEPa-2), on the CAST navigation dataset to produce an 047
013 informed action distribution conditioned on the current ob- 048
014 servation and language instruction. In the second stage, 049
015 we use this policy-derived distribution to warm-start Model 050
016 Predictive Path Integral (MPPI) planning over a separately 051
017 trained JEPa world model, which predicts future latent 052
018 states in the embedding space of the same frozen encoder. 053
019 By initializing the MPPI sampling distribution from the pol- 054
020 icy prior rather than from an uninformed Gaussian, our 055
021 planner converges faster to high-quality action sequences 056
022 that reach the goal. We systematically study the effect of the 057
023 vision encoder backbone, comparing DINOv2 and V-JEPa- 058
024 2, across both the policy and world model components. Ex- 059
025 periments on real-world navigation tasks demonstrate that 060
026 PiJEPa significantly outperforms both standalone policy 061
027 execution and uninformed world model planning, achieving 062
028 improved goal-reaching accuracy and instruction-following 063
029 fidelity. 064

030 1. Introduction

031 Building autonomous agents that can navigate to a goal 065
032 specified by an image and a natural language instruction is a 066
033 long-standing challenge in robotics and computer vision [1– 067
034 3]. Given a current egocentric observation o_t , a goal image 068
035 o_g , and an instruction ℓ (e.g. “move towards the stairs”), 069

the agent must produce a sequence of actions that reliably 036
reaches the goal while respecting the semantics of the in- 037
struction. 038

Recent vision-language-action (VLA) models [3–6] 039
have made remarkable strides toward this objective. Gener- 040
alist policies such as Octo [4] and NoMaD [7] can be fine- 041
tuned on domain-specific navigation data to achieve impres- 042
sive short-horizon control, while the CAST framework [6] 043
further strengthens instruction following through counter- 044
factual language-action augmentation. Despite their effec- 045
tiveness, these reactive policies predict actions in a single 046
forward pass and thus lack the capacity to reason about the 047
long-term consequences of their decisions. 048

A complementary line of work addresses this limitation 049
by learning world models that predict future states condi- 050
tioned on candidate actions, enabling explicit planning to- 051
ward the goal [8–12]. Among these, Joint-Embedding Pre- 052
dictive Architecture World Models (JEPa-WMs) [12–14] 053
are especially attractive: by learning dynamics in the latent 054
space of a frozen pretrained encoder, they support efficient 055
rollouts without pixel-level reconstruction. At inference 056
time, a sampling-based optimizer such as MPPI [15, 16] 057
or CEM [17] searches over action sequences by unrolling 058
the learned predictor and minimizing the embedding-space 059
distance to the encoded goal. However, this search must 060
typically begin from an uninformed prior over a high- 061
dimensional action space, resulting in slow convergence and 062
susceptibility to local minima [13]. 063

We propose PiJEPa, a framework that bridges these two 064
paradigms by using a learned policy to warm-start world 065
model planning. Our key insight is that a finetuned VLA 066
policy, while insufficient for long-horizon reasoning on its 067
own, provides a highly informative *action prior*: it con- 068
centrates probability mass in the region of action space 069
most likely to make progress toward the goal. Initializing 070
the MPPI sampling distribution from this prior allows the 071
world model planner to devote its search budget to refining 072
already-promising trajectories rather than exploring the full 073
action space from scratch. 074

Concretely, our pipeline operates as follows (see Fig- 075

076	ure 1):	instruction-grounded action proposals of a reactive VLA	127
077	1. Policy prior. Given the current observation o_t and in-	policy but uses them as a warm start for a separate plan-	128
078	struction ℓ , a finetuned Octo model produces N_π ac-	ning stage, thereby combining language understanding with	129
079	tion chunk samples via its diffusion head. We transform	long-horizon deliberation.	130
080	these into the world model’s local-frame action space		
081	and compute their empirical mean μ_π and standard devi-		
082	ation σ_π .		
083	2. MPPI planning. A JEPa world model, trained with the	2.2. Latent World Models for Planning	131
084	same frozen encoder, predicts future latent states. We	World models offer a principled route to long-horizon rea-	132
085	run MPPI initialized at (μ_π, σ_π) and optimize against the	soning by enabling agents to simulate the consequences of	133
086	embedding-space distance to the encoded goal o_g .	candidate actions before committing to them. Early ap-	134
087	3. Execution. The first action of the optimized sequence is	proaches such as Navigation World Models (NWM) [11]	135
088	executed, and the process repeats at the next replanning	synthesize pixel-level futures, but predicting dynamics in a	136
089	step.	learned latent space drastically reduces computational cost.	137
090	Our contributions are:	Joint-Embedding Predictive Architecture (JEPa) world	138
091	1. We propose PiJEPa, a unified framework for language-	models built on frozen foundation encoders like DINOv2	139
092	conditioned visual navigation that combines a finetuned	have proven particularly effective for physical planning	140
093	Octo policy with MPPI-based planning over a JEPa	tasks [12, 13]. This family of models has been further ad-	141
094	world model, using the policy output to warm-start the	vanced through continuous latent action formulations [22],	142
095	planner.	principled regularization strategies such as LeJEPa [23],	143
096	2. We train both components with modern pretrained vi-	and joint VLA-latent pretraining [24]. Complementary	144
097	sion encoders (DINOv2 and V-JEPa-2), providing a sys-	efforts have also targeted the efficiency of the planning	145
098	tematic comparison of image-based vs. video-based rep-	loop itself; for instance, One-Step World Model [25] re-	146
099	resentations for policy learning and world modeling in	duces the computational latency of multi-step rollouts, and	147
100	navigation.	benchmarks like Target-Bench [26] evaluate world-model-	148
101	3. We demonstrate on real-world navigation tasks that	based path planning toward text-specified targets. Despite	149
102	policy-guided planning significantly outperforms both	these advances, sampling-based planners such as MPPI and	150
103	standalone policy execution and uninformed MPPI plan-	CEM remain sensitive to their initialization: starting from	151
104	ning, improving goal-reaching accuracy and instruction-	an uninformed Gaussian prior over high-dimensional ac-	152
105	following fidelity.	tion spaces leads to slow convergence and frequent entrap-	153
106		ment in local minima [13]. PiJEPa directly addresses this	154
107	2. Related Work	inference-time bottleneck. Rather than conditioning the	155
108	Our work draws on two active research threads—language-	world model itself on language—which would conflate rep-	156
109	conditioned policies and latent world models—and unifies	resentation learning with instruction grounding—we keep	157
110	them through policy-guided planning. We review each in	the latent dynamics model language-agnostic and instead	158
111	turn.	anchor the MPPI optimization with an instruction-aware ac-	159
112		tion prior drawn from the VLA policy. This decoupled de-	160
113	2.1. Language-Conditioned Navigation Policies	sign lets each component focus on what it does best: the	161
114	The integration of natural language with embodied percep-	policy provides efficient, semantically informed proposals,	162
115	tion has given rise to Vision-Language-Action (VLA) mod-	while the world model evaluates and refines them over ex-	163
116	els that map instructions and visual observations directly to	tended horizons.	164
117	motor commands. Generalist architectures such as Octo [4]		
118	learn broadly from diverse robot data, while recent scalable	3. Method	165
119	designs including Dita [18], LAPA [19], and CLAP [20]	We present PiJEPa, a two-stage approach for instruction-	166
120	push the frontier of reactive instruction-following through	conditioned visual goal navigation. An overview is shown	167
121	latent action representations. In the navigation domain	in Figure 1. Our pipeline consists of three learned com-	168
122	specifically, UrbanNav [21] extends language-guided poli-	ponents: a frozen pretrained vision encoder, an Octo-based	169
123	cies to complex urban environments via web-scale trajec-	policy that generates an informed action prior, and a JEPa	170
124	tory learning. Although these systems achieve strong short-	world model over which we perform MPPI planning warm-	171
125	horizon performance, they share a common limitation: ac-	started by the policy prior.	172
126	tions are selected greedily without explicit reasoning about		
	future states, which degrades performance on tasks that re-	3.1. Problem Formulation	173
	quire multi-step spatial planning. PiJEPa retains the rapid,	At each timestep t , the agent observes an egocentric RGB	174
		image o_t , is given a goal image o_g and a natural language	175

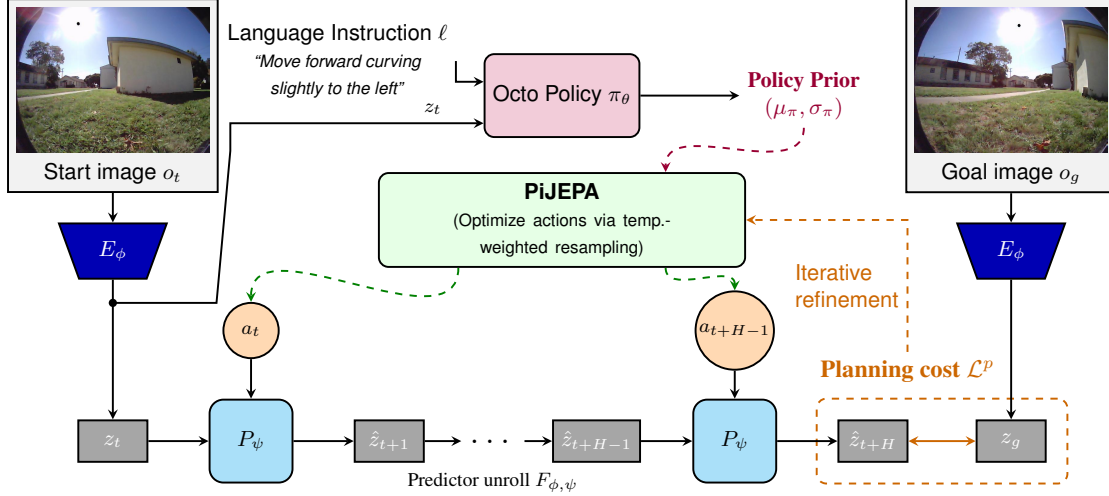


Figure 1. **Overview of PiJEPA.** **Top:** The Octo policy, finetuned with a frozen vision encoder (E_ϕ), takes the current latent observation z_t and instruction ℓ as input and produces action chunk samples via its diffusion head. These are transformed from the global frame to the world model’s local body frame. **Middle:** The policy’s statistics (μ_π, σ_π) warm-start MPPI, which iteratively optimizes the action distribution over J iterations. **Bottom:** The JEPA world model predictor (P_ψ), trained with the same frozen encoder, autoregressively predicts future latent states. The MPPI candidates are scored by unrolling the world model and evaluating the latent-space distance to the encoded goal z_g (Algorithm 1).

176 instruction ℓ , and must produce a sequence of navigation actions
177 $a_{1:H} \in \mathbb{R}^{H \times A}$ over a planning horizon H that guides
178 the robot towards the goal while respecting the instruction
179 semantics.

180 3.2. Policy Prior from Octo

181 The first stage of PiJEPA extracts an action prior from a
182 finetuned Octo policy [4]. Octo is a transformer-based gen-
183 eralist policy with a diffusion action head that models the
184 distribution $\pi_\theta(a|o_t, \ell)$ over action chunks.

185 Both the policy and the world model operate on top of
186 a shared frozen pretrained vision encoder E_ϕ that maps ob-
187 servations to latent tokens $z_t = E_\phi(o_t) \in \mathbb{R}^{n \times d}$, where
188 n is the number of spatial tokens and d the embedding di-
189 mension. The encoder weights ϕ remain frozen throughout,
190 ensuring both components share a consistent representation
191 space.

192 3.3. JEPA World Model

193 The second component is a Joint-Embedding Predictive
194 World Model (JEPA-WM) [12, 13] that predicts future la-
195 tent states given the current state and actions. The world
196 model pairs the same frozen encoder E_ϕ with a learnable
197 predictor P_ψ and action encoder A_ψ . Given a context of en-
198 coded observations and actions, the predictor forecasts the
199 next latent state:

$$200 \hat{z}_{t+1} = P_\psi(z_{t-w:t}, A_\psi(a_{t-w:t})). \quad (1)$$

201 The predictor uses a causal attention mask so it can pre-
202 dict from all context lengths up to a maximum window w .

203 Actions are injected into the predictor via a conditioning
204 mechanism at every layer to prevent vanishing action sig-
205 nals through depth [13].

206 The predictor is trained to minimize the MSE between
207 predicted and target embeddings. Following [13], we use a
208 multi-step rollout loss to improve long-horizon accuracy:

$$209 \mathcal{L}_{\text{total}} = \sum_{k=1}^{K_{\text{roll}}} \frac{1}{B} \sum_{b=1}^B \|F_{\phi,\psi}(z_t^b, a_{t:t+k-1}^b) - z_{t+k}^b\|_2^2, \quad (2)$$

210 where $F_{\phi,\psi}$ denotes the autoregressive unrolling of the pre-
211 predictor (feeding each predicted \hat{z}_{t+j} as context for the next
212 step), trained with truncated backpropagation through time.

213 3.4. Policy-Guided MPPI Planning

214 At test time, we combine the Octo policy prior with the
215 JEPA world model through MPPI-based planning [15, 16].
216 We first draw N_π action chunk samples from the policy
217 diffusion head. Because the policy and world model may
218 operate in different action coordinate frames, each sam-
219 ple is transformed by a mapping T into the world model
220 frame (details in Section 4.1). From the transformed sam-
221 ples $\tilde{a}^{(i)} = T(a^{(i)})$, we compute the initial MPPI mean and
222 standard deviation:

$$223 \mu^0 = \frac{1}{N_\pi} \sum_{i=1}^{N_\pi} \tilde{a}^{(i)}, \quad \sigma^0 = \text{clamp}\left(\text{std}(\{\tilde{a}^{(i)}\}), \sigma_{\min}, \sigma_{\max}\right), \quad (3)$$

224 which defines a policy-informed initialization over action
225 sequences.

226 **Planning objective.** Given the encoded current observation
227 z_t and goal encoding z_g , we evaluate a candidate action se-
228 quence $a_{1:H}$ by the terminal latent distance

$$229 \quad \mathcal{L}^p(z_t, a_{1:H}, z_g) = \frac{1}{n} \sum_{i=1}^n \left\| \hat{z}_{t+H}^{(i)} - z_g^{(i)} \right\|_2^2, \quad (4)$$

230 where $\hat{z}_{t+H} = F_{\phi, \psi}(z_t, a_{1:H})$ is obtained by autoregres-
231 sively rolling out the world model for H steps.

232 **MPPI refinement.** Starting from (μ^0, σ^0) , MPPI itera-
233 tively samples candidate action sequences, evaluates them
234 using Eq. 4, and updates the sampling distribution using
235 temperature-weighted elite trajectories. At iteration j , elite
236 weights are computed as

$$237 \quad w_k = \frac{\exp(\lambda(c_{\min} - c_k))}{\sum_{k'=1}^K \exp(\lambda(c_{\min} - c_{k'}))}, \quad (5)$$

238 and the Gaussian parameters are updated by

$$\mu^{j+1} = \sum_k w_k a^{(k)}, \quad (\sigma^{j+1})^2 = \sum_k w_k (a^{(k)} - \mu^{j+1})^2,$$

239 with σ^{j+1} clamped to $[\sigma_{\min}, \sigma_{\max}]$. After J iterations, one
240 elite trajectory is sampled proportionally to w_k and exe-
241 cuted. The full procedure is summarized in Algorithm 1.

242 **Conditioning.** The language instruction ℓ enters only
243 through the Octo policy, which shapes the initial planning
244 distribution. The world model is language-agnostic and
245 models only visual-action dynamics. The goal image o_g
246 is encoded into z_g and enters through the planning cost in
247 Eq. 4.
248

249 4. Experiments

250 4.1. Experimental Setup

251 **Dataset.** We train and evaluate on the CAST dataset [6], a
252 large-scale visual navigation dataset augmented with coun-
253 terfactual instruction-action pairs. CAST addresses the
254 problem of posterior collapse—where the policy ignores the
255 language instruction because the observation alone suffices
256 to predict the action—by using a VLM to generate alter-
257 native instructions and an atomic policy to produce corre-
258 sponding counterfactual action labels. Each action is a 4-
259 dimensional vector $a = (\Delta x, \Delta y, \sin \Delta \phi, \cos \Delta \phi)$ spec-
260 ifying the robot’s displacement and heading change; the
261 $(\sin \Delta \phi, \cos \Delta \phi)$ encoding avoids angle discontinuities at
262 $\pm \pi$. Actions are normalized to $[-1, 1]$ using bounds nor-
263 malization based on the 1st and 99th percentile statistics.

264 **Vision encoders.** We study two frozen pretrained encoder
265 families: *DINOv2 ViT-S* [27], a self-supervised image en-
266 coder with strong spatial and object segmentation features;
267 and *V-JEPA-2 ViT-L* [28], a self-supervised video encoder
268 that captures temporal dynamics. For V-JEPA-2, we follow

Algorithm 1 PiJEPa: Policy-Guided World Model MPPI Planning

Require: Obs. o_t , goal o_g , instruction ℓ , policy π_θ , world model P_ψ , encoder E_ϕ

Require: Horizon H , samples N , elites K , iterations J , temperature λ , policy samples N_π

Stage 1: Policy prior

```

1:  $z_t \leftarrow E_\phi(o_t), \quad z_g \leftarrow E_\phi(o_g)$ 
2: for  $i = 1, \dots, N_\pi$  do
3:    $a^{(i)} \sim \pi_\theta(\cdot | o_t, \ell)$  ▷ diffusion sampling
4:    $\tilde{a}^{(i)} \leftarrow T(a^{(i)})$  ▷ coordinate transform
5: end for
6:  $\mu^0 \leftarrow \text{mean}(\{\tilde{a}^{(i)}\})$ 
7:  $\sigma^0 \leftarrow \text{clamp}(\text{std}(\{\tilde{a}^{(i)}\}), \sigma_{\min}, \sigma_{\max})$ 

```

Stage 2: MPPI planning

```

8: for  $j = 0, \dots, J - 1$  do
9:   for  $i = 1, \dots, N$  do
10:     $a_{1:H}^{(i)} \leftarrow \text{clamp}(\mu^j + \sigma^j \odot \epsilon^{(i)}, -1, 1)$ 
11:     $\epsilon^{(i)} \sim \mathcal{N}(0, I)$ 
12:     $\hat{z} \leftarrow z_t$ 
13:    for  $h = 1, \dots, H$  do
14:       $\hat{z} \leftarrow P_\psi(\hat{z}, a_h^{(i)})$  ▷ autoregressive rollout
15:    end for
16:     $c^{(i)} \leftarrow \|\hat{z} - z_g\|_2^2 / n$  ▷ Eq. 4
17:  end for
18:   $\mathcal{E} \leftarrow \text{top-}K \text{ by lowest } c^{(i)}$ 
19:   $w_k \leftarrow \exp(\lambda(c_{\min} - c_k)) / \sum_{k'} \exp(\lambda(c_{\min} - c_{k'}))$ 
20:   $\mu^{j+1} \leftarrow \sum_k w_k a^{(k)}$ 
21:   $\sigma^{j+1} \leftarrow \text{clamp}\left(\sqrt{\sum_k w_k (a^{(k)} - \mu^{j+1})^2}, \sigma_{\min}, \sigma_{\max}\right)$ 
22: end for
23: Sample  $k^* \in \mathcal{E}$  with probability  $w_{k^*}$ 
24: return  $a_{1:H}^{(k^*)}$ 

```

the frame-duplication strategy of [13], encoding each frame as a duplicated 2-frame video. We apply layer normalization to the encoder output to stabilize prediction targets and the planning cost landscape.

Policy training. We finetune the Octo-Small model [4] on CAST, replacing its original visual encoder with the frozen pretrained encoder E_ϕ followed by a learnable projection $W_{\text{proj}} \in \mathbb{R}^{d \times d_{\text{Octo}}}$. Language instructions are encoded with a pretrained T5 model [29] (16 tokens). The diffusion action head is a 3-layer MLP with residual connections, trained with the DDPM objective [30] using a cosine noise schedule. Training uses the CAST-augmented data including both original and counterfactual trajectory segments.

Coordinate transform. The Octo policy outputs global-frame displacements, while the world model expects local body-frame actions. We bridge this gap by accumulating

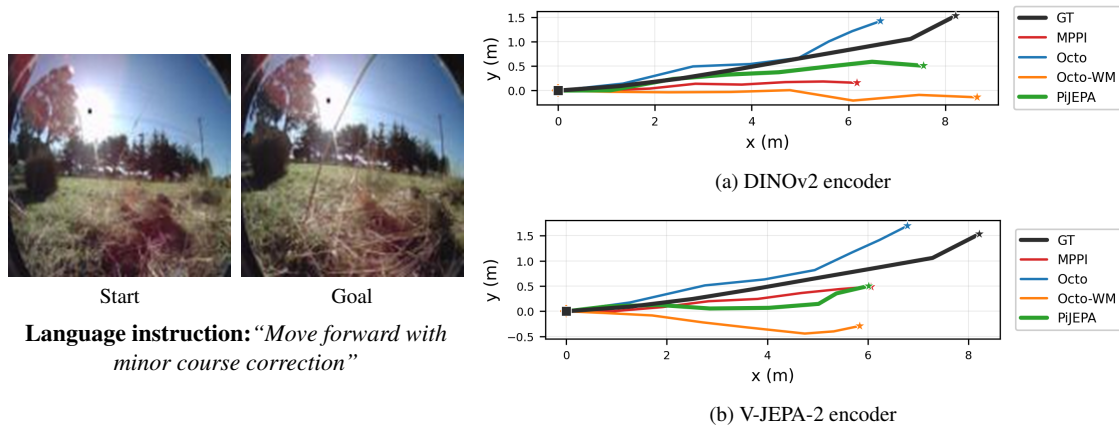


Figure 2. **Qualitative trajectory comparison.** Given a start observation, goal image, and language instruction (left), we compare trajectories produced by each method under two encoder backbones (right). The black curve (GT) shows the ground-truth path; colored curves show MPPI (red), Octo policy (blue), Octo-WM scoring (orange), and PiJEPA (green). Stars mark final positions. PiJEPA most closely tracks the ground truth in both settings.

285 headings from the $(\sin \Delta\phi, \cos \Delta\phi)$ components and rotat- 315
 286 ing the displacement into the robot’s local frame at each 316
 287 timestep:

$$\phi_t = \sum_{\tau=0}^{t-1} \text{atan2}(\sin \Delta\phi_\tau, \cos \Delta\phi_\tau),$$

$$\begin{pmatrix} \Delta x_t^\ell \\ \Delta y_t^\ell \end{pmatrix} = \begin{pmatrix} \cos \phi_t & \sin \phi_t \\ -\sin \phi_t & \cos \phi_t \end{pmatrix} \begin{pmatrix} \Delta x_t \\ \Delta y_t \end{pmatrix}. \quad (7)$$

289 The rotated actions are then normalized to $[-1, 1]$ using the 318
 290 dataset bounds statistics. 319

291 **World model training.** The JEPA world model predictor 320
 292 is a ViT with frame-causal attention and Adaptive Layer 321
 293 Normalization (AdaLN) action conditioning [13, 31], which 322
 294 modulates scale and shift at every transformer block. It is 323
 295 trained on CAST trajectory segments with the multi-step 324
 296 rollout loss (Eq. 2, $K_{\text{roll}} = 2$) using truncated backpropaga- 325
 297 tion through time. The encoder weights remain frozen; only 326
 298 the predictor and action encoder parameters are updated. 327

299 **Encoder variants.** We train four model configurations to 328
 300 systematically study encoder choice: (i) DINOv2-Policy + 329
 301 DINOv2-WM, the most natural consistent-space pairing; 330
 302 (ii) V-JEPA-2-Policy + V-JEPA-2-WM, leveraging tempo- 331
 303 ral features in both components; (iii) DINOv2-Policy + V- 332
 304 JEPA-2-WM; and (iv) V-JEPA-2-Policy + DINOv2-WM. 333
 305 The cross-encoder configurations test whether the action 334
 306 prior from one representation space transfers to planning 335
 307 in another. Prior work [13] found DINOv2 excels at fine- 336
 308 grained spatial reasoning while V-JEPA-2 captures richer 337
 309 temporal structure; evaluating all four combinations reveals 338
 310 which properties matter for the policy prior versus latent 339
 311 dynamics. 340

312 **MPPI hyperparameters.** We use $J = 2$ MPPI iterations, 341
 313 $N = 32$ candidate samples, $K = 4$ elites, inverse tem- 342
 314 perature $\lambda = 0.8$, and $N_\pi = 4$ Octo diffusion samples 343
 344

315 for the policy prior. The standard deviation is clamped to 316
 316 $[\sigma_{\min}, \sigma_{\max}] = [0.01, 0.05]$.

4.2. Baselines and Evaluation 317

To isolate the contribution of each component, we compare 318
 four planning strategies: 319

- **Default MPPI:** Standard MPPI initialized with $\mu^0 = \mathbf{0}$, 320
 $\sigma^0 = \sigma_{\max} \mathbf{1}$, serving as a pure world model planning 321
 baseline. 322
- **PiJEPA (ours):** Warm-started MPPI with $\mu^0 = \mu_\pi$, $\sigma^0 = 323$
 σ_π from the policy prior (Algorithm 1). 324
- **Octo-WM:** Draws N_π samples from the policy, evaluates 325
 each via a full world model rollout, and selects the sam- 326
 ple with the lowest cost \mathcal{L}^p —using the world model for 327
 scoring without trajectory optimization. 328
- **Octo:** Octo predicted actions with no world model in- 329
 volvement, serving as a pure reactive policy baseline. 330

We report Absolute Trajectory Error (ATE) and Rela- 331
 tive Pose Error (RPE) for both XY position (in meters) 332
 and heading (in degrees), computed between predicted and 333
 ground-truth trajectories on the CAST validation set. 334

4.3. Results 335

Tables 1 and 2 report the full set of trajectory metrics for 336
 the DINOv2 and V-JEPA-2 encoder configurations, respec- 337
 tively. We discuss the main findings below. 338

PiJEPA achieves the best positional accuracy across 339
both encoders. PiJEPA leads on ATE XY metrics in both 340
 configurations: 1.78 m RMSE and 3.12 m Final with DI- 341
 NOv2, and 1.65 m RMSE and 1.32 m Mean with V-JEPA- 342
 2—the lowest across all settings. Octo-WM narrowly edges 343
 PiJEPA on V-JEPA-2 Final position (2.87 m vs. 2.88 m), but 344

Table 1. **DINOv2 ViT-S** on CAST validation episodes. Best in **bold**. All ↓.

	ATE XY (m)			ATE Hdg (°)			RPE XY (m)		RPE Hdg (°)	
	RMSE	Mean	Final	RMSE	Mean	Final	RMSE	Mean	RMSE	Mean
Octo	1.98	1.65	3.19	20.08	16.76	28.55	0.61	0.57	8.29	6.94
MPPI	1.85	1.48	3.23	18.79	15.68	27.86	0.55	0.51	5.70	5.11
Octo-WM	1.80	1.43	3.15	22.16	18.60	29.16	0.56	0.51	7.51	6.45
PiJEPa	1.78	1.42	3.12	19.89	16.52	28.63	0.56	0.51	7.05	6.07

Table 2. **V-JEPA-2 ViT-L** on CAST validation episodes. Best in **bold**. All ↓.

	ATE XY (m)			ATE Hdg (°)			RPE XY (m)		RPE Hdg (°)	
	RMSE	Mean	Final	RMSE	Mean	Final	RMSE	Mean	RMSE	Mean
Octo	1.72	1.36	3.02	18.72	15.64	27.72	0.53	0.48	6.086	5.34
MPPI	1.87	1.50	3.29	19.15	15.98	28.50	0.55	0.52	5.75	5.15
Octo-WM	1.67	1.35	2.87	20.63	17.27	28.06	0.54	0.49	6.93	5.99
PiJEPa	1.65	1.32	2.88	19.13	15.95	27.93	0.53	0.49	6.46	5.62

345 PiJEPa dominates on the remaining trajectory-level met-
 346 rics, confirming that warm-starting MPPI with a policy prior
 347 enables the planner to refine already-promising trajectories
 348 rather than searching from scratch.

349 **Uninformed MPPI struggles with position but excels at**
 350 **heading.** MPPI exhibits a striking divergence across met-
 351 ric types: it is the weakest method on ATE XY under
 352 V-JEPA-2 (RMSE 1.87 m, Final 3.29 m), confirming that
 353 uninformed initialization wastes the sampling budget, yet
 354 it achieves the best heading metrics across both encoders
 355 (e.g., DINOv2: ATE Hdg RMSE 18.79°, RPE Hdg Mean
 356 5.11°). We attribute this to the embedding-space objective
 357 implicitly regularizing heading alignment even when trans-
 358 lational accuracy is poor.

359 **World model scoring captures much of the planning**
 360 **benefit.** Octo-WM provides large positional gains over
 361 the raw Octo policy—with V-JEPA-2, ATE XY Final drops
 362 from 3.02 m to 2.87 m—suggesting that even without itera-
 363 tive optimization, using the world model to filter poor policy
 364 proposals is a powerful strategy. The relatively small gap
 365 between Octo-WM and PiJEPa indicates that filtering ac-
 366 counts for much of the benefit, though the additional MPPI
 367 refinement still yields the best overall trajectory accuracy.

368 **The reactive policy provides strong local control.** De-
 369 spite lacking look-ahead, Octo remains competitive on step-
 370 level metrics: with V-JEPA-2 it achieves the best RPE XY
 371 Mean (0.48 m) and strongest ATE Heading (RMSE 18.72°,
 372 Final 27.72°). Planning thus contributes its main benefit at
 373 the trajectory level, where cumulative errors compound.

V-JEPA-2 yields stronger overall performance. V-
 JEPA-2 with PiJEPa achieves the best overall positional ac-
 curacy, and even its Octo baseline (ATE XY RMSE 1.72 m)
 surpasses the best DINOv2 method (1.78 m). However,
 uninformed MPPI performs worst under V-JEPA-2 (Final
 3.29 m vs. 3.23 m for DINOv2), suggesting its richer la-
 tent dynamics create a harder optimization landscape for
 uninformed search. Once anchored by the policy prior, this
 space becomes an asset: the gap between PiJEPa and un-
 informed MPPI on ATE XY Final grows from 0.11 m (DI-
 NOv2) to 0.41 m (V-JEPA-2).

Planning latency. The Octo policy requires approxi-
 mately 2.13 s to generate its action proposals via diffusion
 sampling, while the MPPI planning stage adds only 0.35 s,
 bringing the total PiJEPa inference time to roughly 2.48 s
 for an 8-step trajectory. The lightweight MPPI overhead
 demonstrates that the planning stage is practical to layer
 on top of the policy, with the vast majority of latency at-
 tributable to the diffusion-based action sampling rather than
 the world model rollouts.

Qualitative results and failure analysis. Figure 2 shows
 that PiJEPa (green) most closely tracks the ground-truth
 path under both encoders, while MPPI (red) deviates sub-
 stantially and Octo (blue) drifts over longer horizons. Fig-
 ure 3 presents a failure case with the ambiguous instruc-
 tion “*Follow the building,*” where the Octo policy veers
 toward the wrong building. The world model itself gets
 stuck in its rollouts, and the predicted latent states remain
 nearly identical across the planning horizon regardless of
 the sampled actions, leaving the planner unable to make
 progress. PiJEPa partially mitigates this through the policy

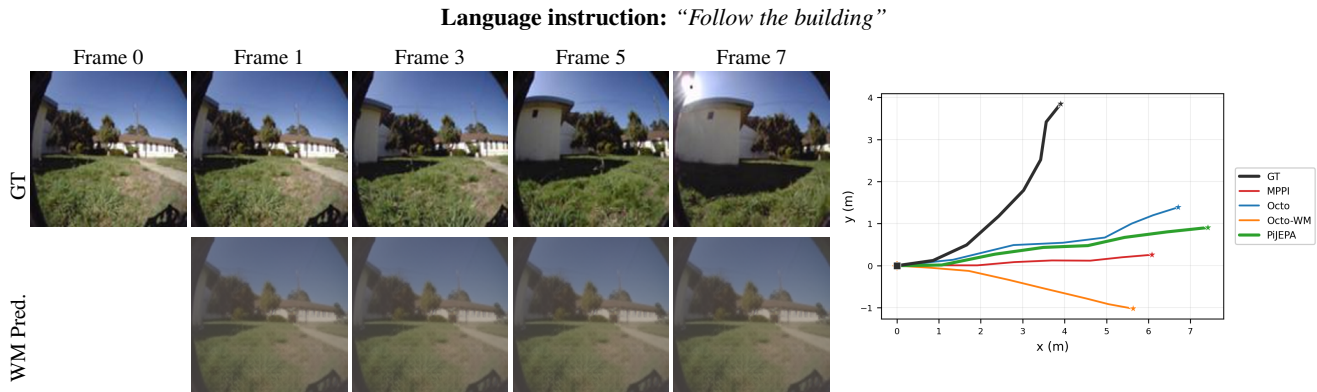


Figure 3. **Failure case analysis.** The language instruction “*Follow the building*” is inherently ambiguous, as multiple buildings are visible in the scene. The Octo policy (blue) misinterprets the referent and veers toward a different building, illustrating how vague instructions can mislead reactive policies that lack long-horizon reasoning. Meanwhile, the world model fails to make meaningful progress because it becomes stuck in its rollouts, predicting nearly identical latent states, which causes the planner to stagnate. The WM Pred. row confirms this directly. The predicted observations remain largely unchanged across the planning horizon. PiJEPa (green) partially mitigates both issues by grounding the planner with a policy-derived prior, though it still undershoots the ground-truth trajectory.

405 prior, which bypasses the stalled rollouts by anchoring the
 406 search in a productive region of action space, though it still
 407 undershoots the ground truth, indicating that ambiguous in-
 408 structions and world model stagnation remain challenging.

409 5. Conclusion

410 We have presented PiJEPa, a framework for instruction-
 411 conditioned visual navigation that warm-starts MPPI plan-
 412 ning over a JEPa world model using a finetuned VLA
 413 policy prior. PiJEPa achieves the best positional accu-
 414 racy across both DINOv2 and V-JEPa-2 encoders, outper-
 415 forming reactive policies, uninformed planning, and world
 416 model scoring baselines. Our analysis reveals a natural di-
 417 vision of labor: the reactive policy excels at local control
 418 and heading prediction, uninformed MPPI achieves strong
 419 heading alignment through its embedding-space objective,
 420 and world model scoring captures much of the look-ahead
 421 benefit by filtering poor proposals—PiJEPa combines these
 422 strengths, with the additional MPPI refinement yielding the
 423 best trajectory-level coherence. The MPPI planning latency
 424 overhead is negligible relative to the accuracy gains.

425 However, our failure analysis reveals that the world
 426 model can become stuck in its rollouts, producing static
 427 latent predictions that render the planner ineffective. Ad-
 428 dressing this stagnation through improved dynamics archi-
 429 tectures or diversity-promoting rollout mechanisms is a key
 430 direction for future work, alongside intermediate waypoint
 431 costs, richer action spaces, and closed-loop replanning.

432 References

433 [1] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-
 434 Fei, and A. Farhadi, “Target-driven visual navigation in in-

door scenes using deep reinforcement learning,” in *IEEE In-
 ternational Conference on Robotics and Automation*, 2017. 435
 436

- 437 1 437
- [2] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, 438
 N. Sünderhauf, I. Reid, S. Gould, and A. van den Hen- 439
 gel, “Vision-and-language navigation: Interpreting visually- 440
 grounded navigation instructions in real environments,” in 441
IEEE Conference on Computer Vision and Pattern Recogni- 442
tion, 2018. 443
- [3] D. Shah, A. Sridhar, N. Dashora, K. Stachowicz, K. Black, 444
 N. Hirose, and S. Levine, “ViNT: A foundation model for 445
 visual navigation,” in *Conference on Robot Learning*, 2023. 446
 1 447
- [4] O. M. Team, D. Ghosh, H. R. Walke, K. Pertsch, K. Black, 448
 O. Mees, S. Dasari, J. Hejna, T. Kreiman, C. Xu, J. Luo, Y. L. 449
 Tan, P. R. Sanketi, Q. Vuong, T. Xiao, D. Sadigh, C. Finn, 450
 and S. Levine, “Octo: An open-source generalist robot pol- 451
 icy,” *ArXiv*, vol. abs/2405.12213, 2024. 1, 2, 3, 4 452
- [5] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, 453
 K. Choromanski, *et al.*, “RT-2: Vision-language-action mod- 454
 els transfer web knowledge to robotic control,” in *Confer- 455*
ence on Robot Learning, 2023. 456
- [6] C. Glossop, W. Chen, A. Bhorkar, D. Shah, and S. Levine, 457
 “CAST: Counterfactual labels improve instruction fol- 458
 lowing in vision-language-action models,” *arXiv preprint 459*
arXiv:2508.13446, 2025. 1, 4 460
- [7] A. Sridhar, D. Shah, C. Glossop, and S. Levine, “NoMaD: 461
 Goal masked diffusion policies for navigation and explo- 462
 ration,” *arXiv preprint arXiv:2310.07896*, 2023. 1 463
- [8] Y. LeCun, “A path towards autonomous machine intelli- 464
 gence,” *OpenReview preprint*, 2022. 1 465
- [9] D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap, “Master- 466
 ing diverse domains through world models,” in *International 467*
Conference on Machine Learning, 2024. 468

- 469 [10] N. Hansen, H. Su, and X. Wang, “TD-MPC2: Scalable, robust world models for continuous control,” *arXiv preprint arXiv:2310.16828*, 2024. 525
- 470 526
- 471 527
- 472 [11] A. Bar, G. Zhou, D. Tran, T. Darrell, and Y. LeCun, “Navigation world models,” *arXiv preprint arXiv:2412.03572*, 2024. 528
- 473 529
- 474 2 530
- 475 [12] G. Zhou, H. Pan, Y. LeCun, and L. Pinto, “Dino-wm: World models on pre-trained visual features enable zero-shot planning,” 2024. 1, 2, 3 531
- 476 532
- 477 533
- 478 [13] B. Terver, T.-Y. Yang, J. Ponce, A. Bardes, and Y. LeCun, “What drives success in physical planning with joint-embedding predictive world models?,” 2025. 1, 2, 3, 4, 5 534
- 479 535
- 480 536
- 481 [14] M. Assran *et al.*, “V-JEPA-2-AC: Vision joint-embedding predictive architecture with action conditioning,” *arXiv preprint*, 2025. 1 537
- 482 538
- 483 539
- 484 [15] G. Williams, P. Drews, B. Goldfain, J. M. Rehg, and E. A. Theodorou, “Aggressive driving with model predictive path integral control,” in *IEEE International Conference on Robotics and Automation*, 2016. 1, 3 540
- 485 541
- 486 542
- 487 543
- 488 [16] G. Williams, N. Wagener, B. Goldfain, P. Drews, J. M. Rehg, B. Boots, and E. A. Theodorou, “Information theoretic MPC for model-based reinforcement learning,” *IEEE International Conference on Robotics and Automation*, 2017. 1, 3 544
- 489 545
- 490 546
- 491 547
- 492 548
- 493 [17] R. Rubinstein, “The cross-entropy method for combinatorial and continuous optimization,” *Methodology and Computing in Applied Probability*, vol. 1, pp. 127–190, 1999. 1
- 494 549
- 495 550
- 496 [18] Z. Hou, T. Zhang, Y. Xiong, H. Duan, H. Pu, R. Tong, C. Zhao, X. Zhu, Y. Qiao, J. Dai, and Y. Chen, “Dita: Scaling diffusion transformer for generalist vision-language-action policy,” *ArXiv*, vol. abs/2503.19757, 2025. 2 551
- 497 552
- 498 553
- 499 554
- 500 [19] S. Ye, J. Jang, B. Jeon, S. J. Joo, J. Yang, B. Peng, A. Mandekar, R. Tan, Y.-W. Chao, B. Y. Lin, L. Lidén, K. Lee, J. Gao, L. S. Zettlemoyer, D. Fox, and M. Seo, “Latent action pretraining from videos,” *ArXiv*, vol. abs/2410.11758, 2024. 2 555
- 501 556
- 502 557
- 503 558
- 504 559
- 505 [20] C. Zhang, J. Wang, Z. Gao, Y. Su, T. Dai, C. Zhou, J. Lu, and Y. Tang, “Clap: Contrastive latent action pretraining for learning vision-language-action models from human videos,” *ArXiv*, vol. abs/2601.04061, 2026. 2 560
- 506 561
- 507 562
- 508 563
- 509 [21] Y. Mei, Y. Yang, L. Guo, Q. Wang, M. Yu, X. He, W. Wu, and J. Liu, “Urbannav: Learning language-guided urban navigation from web-scale human trajectories,” *ArXiv*, vol. abs/2512.09607, 2025. 2 564
- 510 565
- 511 566
- 512 567
- 513 [22] Q. Garrido, T. Nagarajan, B. Terver, N. Ballas, Y. LeCun, and M. Rabbat, “Learning latent action world models in the wild,” *ArXiv*, vol. abs/2601.05230, 2026. 2 568
- 514 569
- 515 570
- 516 [23] R. Balestriero and Y. LeCun, “Lejepa: Provable and scalable self-supervised learning without the heuristics,” *ArXiv*, vol. abs/2511.08544, 2025. 2 571
- 517 572
- 518 573
- 519 [24] J. Sun, W. Zhang, Z. Qi, S. Ren, Z. Liu, H. Zhu, G. Sun, X. Jin, and Z. Chen, “Vla-jepa: Enhancing vision-language-action model with latent world model,” 2026. 2 574
- 520 575
- 521 576
- 522 [25] W. Shen, Z. Meng, J. Ma, M. Zhou, and D. Xiang, “An efficient and multi-modal navigation system with one-step world model,” *ArXiv*, vol. abs/2601.12277, 2026. 2 577
- 523 578
- 524 579
- [26] D. Wang, H. Ye, Z. Liang, Z. Sun, Z. Lu, Y. Zhang, Y. Zhao, Y. Gao, M. Seeger, F. Schäfer, H. Qin, W. Li, L. Palmieri, F. Jahncke, M. Piccinini, and J. Betz, “Target-bench: Can world models achieve mapless path planning with semantic targets?,” *ArXiv*, vol. abs/2511.17792, 2025. 2 580
- [27] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, *et al.*, “DINOv2: Learning robust visual features without supervision,” *Transactions on Machine Learning Research*, 2024. 4 581
- [28] M. Assran *et al.*, “V-JEPA-2: Self-supervised video models enable understanding, generation and planning,” *arXiv preprint*, 2025. 4 582
- [29] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020. 4 583
- [30] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” in *Advances in Neural Information Processing Systems*, 2020. 4 584
- [31] W. Peebles and S. Xie, “Scalable diffusion models with transformers,” in *International Conference on Computer Vision*, 2023. 5 585