

Learning Latent Dynamics for Planning from Pixels Danijar Hafner¹², Timothy Lillicrap³, Ian Fischer⁴, Ruben Villegas¹⁵, David Ha¹, Honglak Lee¹, James Davidson¹ ¹Google Brain ² University of Toronto ³ DeepMind ⁴ Google Research ⁵ University of Michigan

1 We introduce PlaNet

Planning is a powerful approach when the rules of the environment are known, for example in AlphaGo.



5

To plan in unknown environments, the agent needs to learn a world model over time from experience.



- By planning in latent space, PlaNet reaching top model-free performance on visual control tasks in 200X fewer episodes.
- We also train one agent to make accurate video predictions and solve all six considered tasks without task labels.

4 Recurrent State Space Model

Deterministic states can pass information far into the future, needed for accurate predictions.

Stochastic states can predict multiple futures, needed for partial observability and model uncertainty.

Our experiments show that having both stochastic and deterministic state elements is crucial for planning performance.



Contact: mail@danijar.com Twitter: @danijarh



2 Latent Dynamics Model

Our model predicts future images and rewards using a sequence of compact latent states, trained as a sequential VAE.

Reconstructing images provides a rich training signal but is not needed during planning.



5 Comparison to Model-Free Agents





(a) Cartpole

(b) Reacher

(c) Cheetah

Our visual control tasks include partial observability, contact dynamics, sparse rewards, and many degrees of freedom.

PlaNet reaches the performance of top model-free agents in 200x fewer episodes on average on these tasks.

Method	Modality	Episodes	Cartpole Swing Up	Reacher Easy	Cheetah Run	Finger Spin	Cup Catch	Walker Walk
A3C	proprioceptive	100,000	558	285	214	129	105	311
D4PG	pixels	100,000	862	967	524	985	980	968
PlaNet (ours)	pixels	1,000	821	832	662	700	930	951
Data efficiency gain PlaNet over D4PG (factor)			250	40	500+	300	100	90

Training time: 1 day on a single GPU



(d) Finger

(f) Walker

3 Planning in Latent Space

PlaNet selects actions by simple population based search, without needing a policy network. We replan at every step.

We only need future rewards not images, allowing to evaluate thousands of action sequences in parallel in a large batch.



6 Comparison of Model Designs

Our paper includes further experiments for the planning method, learned representations, multi-step predictions, and more.





The plots show median performance and percentiles 5 to 95 over 5 seeds and 10 episodes each.

Project website with videos: danijar.com/planet