

000 IMPROVING WORLD MODELS USING LINEAR PROBES 001

002
003 **Anonymous authors**

004 Paper under double-blind review
005
006
007

ABSTRACT

010 Developing effective world models is crucial for creating artificial agents that can
011 reason about and navigate complex environments. In this paper, we investigate a
012 technique for encouraging the development of a world model in a network trained
013 end-to-end to predict the next observation. Focusing on a toy environment based
014 on the Flappy Bird game, where the agent receives only LIDAR measurements
015 as observations, we explore the effect of adding a linear probe component to the
016 network’s loss function. This additional term encourages the network to encode a
017 subset of the true underlying world features into its hidden state. Our experiments
018 demonstrate that this supervision technique improves both training and test per-
019 formance, enhances training stability, and results in more easily decodable world
020 features – even for those world features which were not included in the train-
021 ing. Furthermore, we observe reduced distribution drift in networks trained with
022 the linear probe, particularly during high-variability phases of the game (flying
023 between successive pipe encounters). While the scaling laws for model size and
024 training time maintained the same slope with and without the linear probe, the loss
025 curve when training with the linear probe was shifted down, indicating improved
026 performance at a given model size or training duration. Including the world fea-
027 tures loss component roughly corresponded to doubling the model size. This sug-
028 gests that the linear probe technique is particularly beneficial in compute-limited
029 settings or when aiming to achieve the best performance with smaller models.
030 These findings contribute to our understanding of how to develop more robust
031 and sophisticated world models in artificial agents, paving the way for further
032 advancements in this field.

033 1 INTRODUCTION

034
035 Developing effective world models is a crucial aspect of artificial intelligence, as it can enable agents
036 to make accurate predictions of how the world will unfold, as well as how their actions will influence
037 the world, and plan their actions accordingly (Ha & Schmidhuber, 2018). World models allow
038 agents to reason about the environment and its dynamics. In the context of AI, a “world model”
039 typically refers to a learned approximation of the environment’s underlying dynamics, which may
040 include physical laws, object interactions, and agent-environment relationships. By compressing
041 sensory observations into a structured latent space, world models allow agents to reason about their
042 surroundings and adapt to changing conditions. However, learning robust and accurate world models
043 remains a significant challenge, especially in complex and partially observable environments.

044 Recent advancements in machine learning have shown promise in learning world models through
045 end-to-end training of predictive networks (Hafner et al., 2020). These networks are trained to
046 predict the next observation based on the current observation and action, essentially learning a model
047 of the environment’s dynamics. However, it is unclear whether such networks can develop implicit
048 world models that capture the true underlying structure and features of the environment.

049 In this paper, we consider an example environment – Flappy Bird with LIDAR (Figure 1) – and
050 investigate whether an end-to-end trained predictive network can indeed develop implicit world
051 models. Crucially, we explore techniques to encourage the emergence of such models. Specifically,
052 we focus on the following research question: Can an addition of a linear probe to the network’s loss
053 function, which encourages the decoding of true world features, improve the network’s performance
and the quality of the learned world model?

By exploring these questions and techniques, we aim to contribute to the understanding and development of more effective world models in artificial intelligence, ultimately leading to more capable and adaptable agents.

1.1 RELATED WORK

Deep supervision techniques were applied extensively in training vision models

2 METHODS

2.1 ENVIRONMENT AND TASK DESCRIPTION

We use an experimental environment based on the Flappy Bird game¹ (Figure 1A), where the bird must successfully pass through layers of pipes by performing either a no-op action or a flap action, which propels the bird upward (Figure 1C). The agent receives only LIDAR measurements as observations, represented by a 180-dimensional vector indicating the distance to the nearest obstacle for each ray angle (Figure 1B). Separately, the environment also provides the true underlying world features, such as the bird’s rotation, vertical velocity, and position (Figure 1D-F), as well as information about the pipes.

The agent (bird) receives reward every time it successfully passes through a pair of pipes. The episode ends whenever the bird touches any obstacle (pipe, floor, or ceiling). However, in this paper we focus exclusively on the world model aspect of this environment, and so the rewards are not considered.

For the experiments in this paper, we use 10,000 training rollouts of the environment collected by an expert DQN policy (which turns into a random policy with probability 1% every timestep; this randomness was included to encourage off-policy exploration). For testing, we use 2,000 rollouts of the DQN policy and 2,000 rollouts from a random policy (which chooses action "flap" with probability 7.5% every timestep).

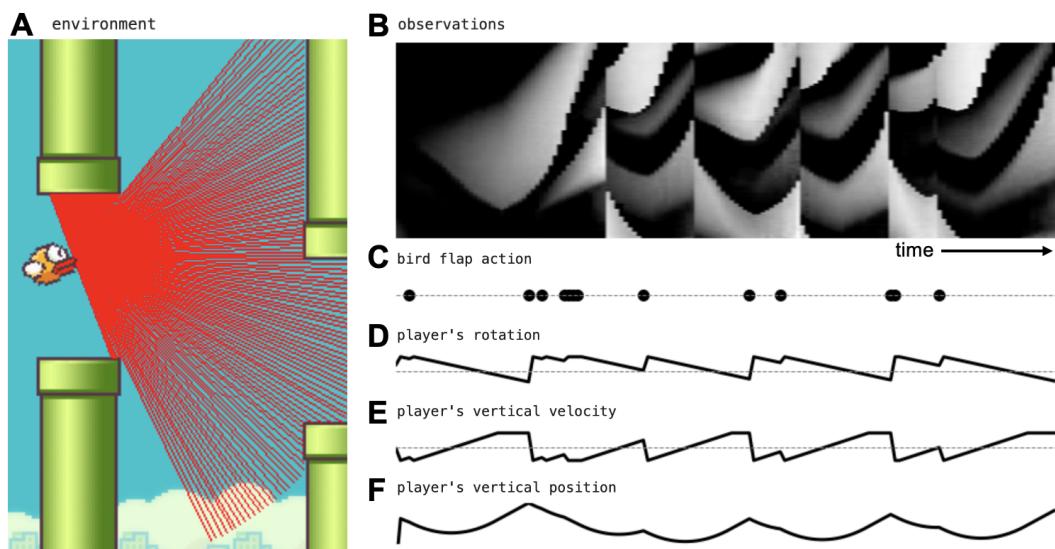


Figure 1: Flappy Bird environment with lidar. (A) The environment. (B) The agent only observes the lidar signal as a function of time. (C) The only available actions are no-op and flap. (D-F) The environment provides true variables of the world, such as the player’s rotation angle, vertical velocity, and position.

¹The original implementation of the environment can be found at <https://github.com/markub3327/flappy-bird-gymnasium>.

108 2.2 NETWORK ARCHITECTURE AND TRAINING
109110 The network architecture is based on that of the world models in the paper by Ha & Schmidhuber
111 (2018).112 First, raw 180-dimensional observations are transformed into an 8-dimensional latent space using
113 an autoencoder (Figure 2A). The autoencoder consists of an encoder multi-layer perceptron (MLP)
114 E that maps the high-dimensional observations to the lower 8-dimensional latent space (outputting
115 the mean and variance of an 8-dimensional factored gaussian distribution), and a decoder MLP D
116 that reconstructs the original observations from the latent representation.117 The autoencoder is trained on the collected rollouts, with L2 regularization, with the following loss:
118

119
$$L_{AE} = \sum_i (\hat{x}_i - x_i)^2 + \frac{1}{2} \sum_i \|E(x_i)\|^2, \quad (1)$$

120

121 where

122
$$\hat{x}_i = D(z_i); \quad z_i \sim \mathcal{N}(\mu_i, \sigma_i^2); \quad E(x_i) = (\mu_i, \sigma_i). \quad (2)$$

123

124 After training, the weights of the autoencoder are frozen.
125126 The world model is a long short-term memory RNN (LSTM; Hochreiter & Schmidhuber (1997);
127 Figure 2B). This RNN takes as input the current latent observation vector (z_t) and action a_t , and
128 outputs the distribution of the next latent vector ($P(\tilde{z}_{t+1})$) and an episode end flag \tilde{e}_{t+1} . The dis-
129 tribution of latent vectors is represented using a mixture of 5 Gaussian variables (following the im-
130 plementation of the mixture density network; Graves (2014); Bishop (1994)) – that is, the network
131 outputs 5 tuples $(\mu^{(i)}, \sigma^{(i)})$ to represent the distribution $P(\tilde{z}_{t+1})$.

132 The loss function is

133
$$L_{pred} = -\frac{1}{n} \sum_t \log \left[\frac{1}{5} \sum_{i=1}^5 P(z_{t+1} | \mu_t^{(i)}, \sigma_t^{(i)}) \right] + \frac{10}{n} \sum_t (e_t - \tilde{e}_t)^2, \quad (3)$$

134

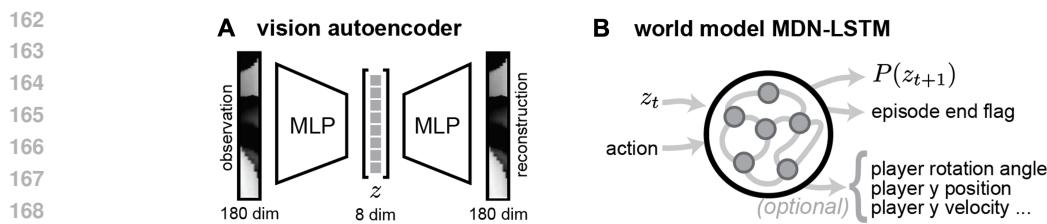
135 where n is the number of timesteps in an episode.
136137 Our key contribution is the (optional) addition of linear probes to the network, which attempt to de-
138 code the true underlying world features from the LSTM’s hidden state. The loss function is modified
139 to include a term proportional to the mean squared error of the linear probes, with a hyperparameter
140 multiplier λ . The total loss is the sum of the predictive loss and the linear probe mean squared error:
141

142
$$L_{total} = L_{pred} + \lambda \cdot MSE(\text{linear probe}) \quad (4)$$

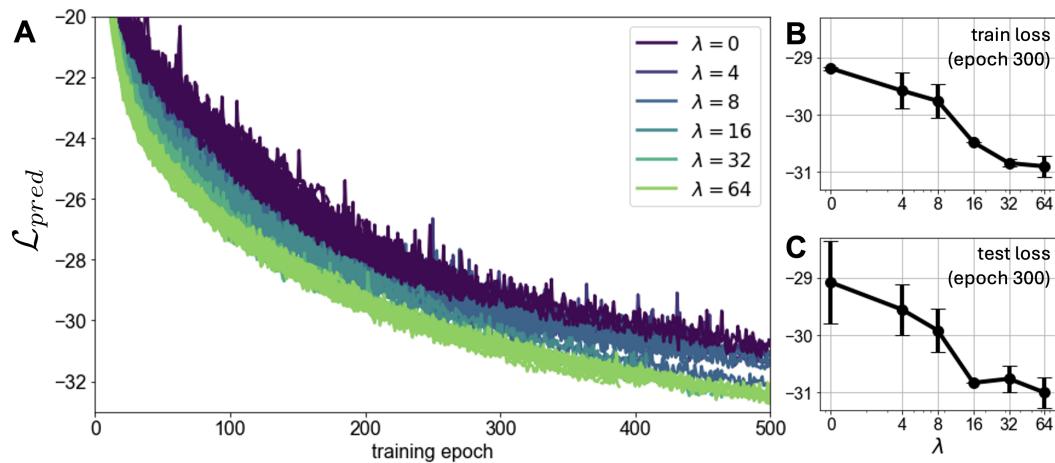
143
$$= L_{pred} + \lambda \sum_{k,t} (W_k r_t - f_{k,t})^2, \quad (5)$$

144

145 where W_k are the linear weights for world feature k ; r_t are the features in the recurrent layer of the
146 LSRM at time t , and $f_{k,t}$ is the value of the world feature k for the timestep t .
147148 3 RESULTS
149150 3.1 IMPROVED NEXT STATE PREDICTION LOSS
151152 To assess the effects of adding linear probes during training, we trained world model RNNs for
153 various values of λ (20 random seeds for each choice of λ). If $\lambda = 0$, the setting is equivalent to
154 training on the predictive loss only, without any linear probe component.
155156 Perhaps paradoxically, increasing the hyperparameter λ , which put more weight on the linear probe
157 component of the loss, did not hurt the predictive loss component. Instead, both the training and
158 test predictive losses decreased as λ was increased (Figure 3A-C). In fact, the predictive loss as a
159 function of λ continued to decrease through values as high as $\lambda = 64$.160 This result suggested that the addition of the linear probe component to the loss function improved
161 the network’s performance without compromising its predictive capabilities. For the rest of this
paper, only networks with $\lambda = 0$ and $\lambda = 64$ will be considered and compared together.



170 Figure 2: Network architecture and training setup. (A) The vision autoencoder compresses the
171 180-dimensional raw observations into an 8-dimensional latent space. (B) The world model MDN-
172 LSTM takes the current latent observation vector and action as inputs, and predicts the distribution
173 of the next latent vector and an episode end flag. Optionally, the LSTM’s hidden state is encouraged
174 to capture true world variables (such as player rotation angle, y position, y velocity, etc.) through a
175 linear probe.



191 Figure 3: The effect of increasing the linear probe weight λ on the original (next latent state pre-
192 diction) loss. Both training (A, B) and test (C) predictive losses decrease as λ increases. For every
193 choice of λ , 20 RNNs were initialized from different random seeds. Error bars in panels (B) and (C)
194 indicate the s.d.

197 3.2 IMPROVED DECODABILITY OF WORLD FEATURES

199 To assess the sophistication of the world model, we examined the performance of linear probes in
200 decoding true world variables from the network’s hidden state.

201 As expected, the world features that were included in the linear probe loss during training are more
202 easily decodable for $\lambda = 64$ than for $\lambda = 0$ (Figure 4; last three panels – player’s vertical position,
203 velocity, and rotation angle).

204 Interestingly however, only for $\lambda = 64$, even features which were not explicitly part of the loss
205 function exhibited decodability above that of an untrained randomly-initialized network (Figure 4;
206 all panels). This finding indicated that the addition of the linear probe component encouraged the
207 network to develop a more comprehensive representation of the underlying world features, perhaps
208 developing a more sophisticated world model.

210 3.3 REDUCED DISTRIBUTION DRIFT

212 Next, we tested the network’s ability to predict how the latent state changes over time, without
213 receiving constant feedback input from the environment. Specifically, we chose a certain timestep,
214 after which the network was cut off from the environment (stopped receiving the true input latent
215 vector every timestep). Instead, after the cutoff timestep, we sampled the next predicted latent state
from the network’s output and fed it back as its input in the subsequent step.

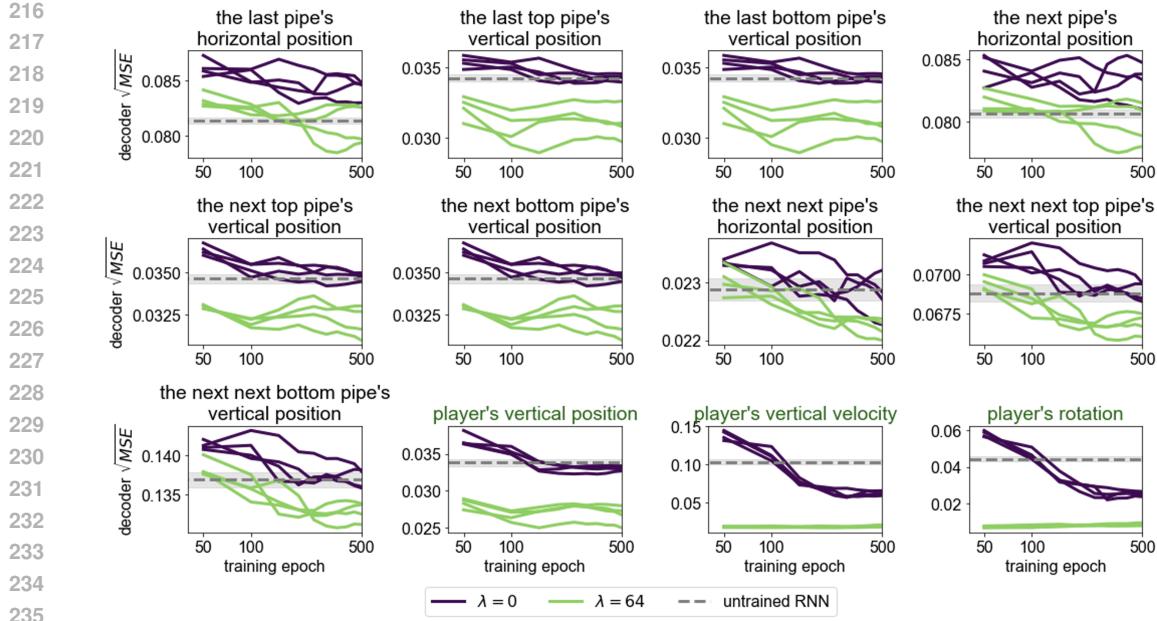


Figure 4: Decodability of world features from the network’s hidden state for $\lambda = 0$ and $\lambda = 64$. Every panel represents one world feature. The last three panels (player’s vertical position, velocity, and rotation; green) represent features which were explicitly included in the loss function for $\lambda = 64$. Note that even the features which were not explicitly part of the loss function exhibit decodability above that of an untrained randomly-initialized network, but only for $\lambda = 64$. For every choice of λ , four RNNs initialized from different random seeds are shown.

Figure 5 compares the distribution drift between networks trained with and without the linear probe. Networks trained with the linear probe ($\lambda = 64$) exhibited smaller growth in the loss of predicting the real latent state, indicating reduced distribution drift, but only for timesteps when the bird was flying between successive pipe encounters ($t=10$ and $t=60$; Figure 5A, C). Interestingly, during the timesteps when the bird was flying through a pair of pipes, this result did not hold. This may have to do with the fact that the changes in game state are much less predictable in those timesteps.

The results were qualitatively equivalent for both the training policy (DQN) and a random policy (Figure 5D-F). This finding suggests effectiveness of the linear probes in maintaining a more accurate prediction of the latent observation state over time.

3.4 SCALING PROPERTIES

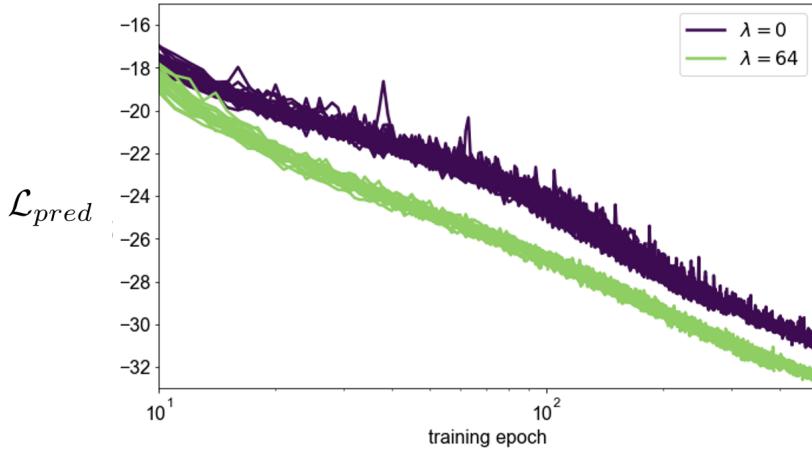
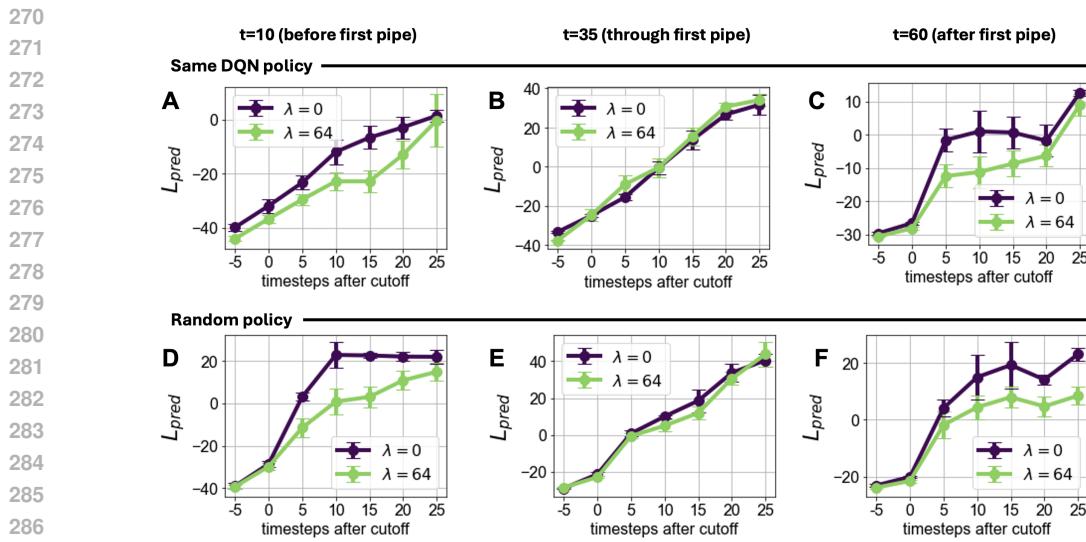
The predictive component of the loss followed a scaling law with respect to training time (Figure 6), for both $\lambda = 0$ and $\lambda = 64$. However, the scaling law for $\lambda = 64$ was smoother, as well as shifted down, compared to the one for networks trained with $\lambda = 0$.

The loss of networks with and without the linear probe also followed a scaling law with respect to model size (Figure 7). However, the curve for networks trained with $\lambda = 64$ was shifted down, such that a network of a given size trained with the linear probe achieves performance roughly equivalent to a network twice the size trained without the probe. This result held true at both epoch 250 and 500.

This finding suggested that the linear probe component provides a consistent performance benefit across different stages of training and different model sizes.

3.5 TRAINING STABILITY

Finally, adding the linear probe mean squared error component to the loss function improved training stability (Supplementary Materials; Figure 8). It resulted in fewer instances of exploding gradients



and training divergences, with the divergences additionally happening later in training for $\lambda = 64$ than for $\lambda = 0$. Networks trained with the linear probe are more likely to reach epoch 500 without diverging (30% vs 7.5% for RNNs of size 128, 62.5% vs 10% for RNNs of size 256; Figure 8A-B, right).

4 DISCUSSION

In this study, we sought to encourage the development of world models in a recurrent neural network trained end-to-end to predict the next observation of an agent. Specifically, we considered adding a

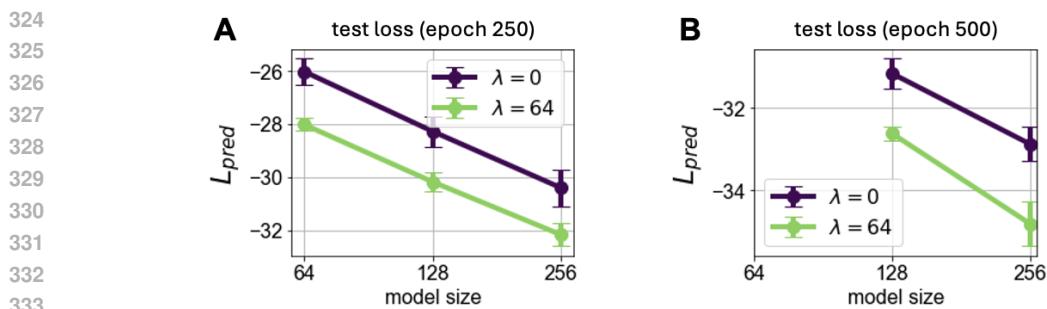


Figure 7: Scaling laws for model size. The performance of networks with and without the linear probe follows a scaling law, but the curve for networks trained with $\lambda = 64$ is shifted down, indicating better performance for a given model size. This holds true at both (A) epoch 250 and (B) epoch 500.

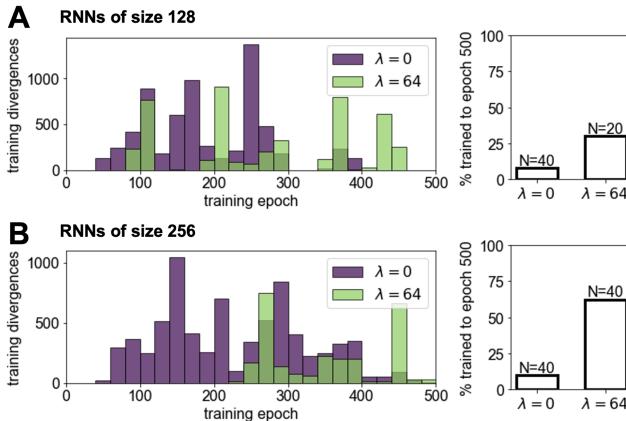


Figure 8: Training stability with and without the linear probe. Networks trained with the linear probe are more likely to reach epoch 500 without diverging. Here, results are shown for RNNs of hidden size 128 (A) and 256 (B).

term to the loss function of the network, which was proportional to the mean squared error of linear probes trying to decode the true world features from the hidden state of the model.

In our experiments, this addition demonstrated improved training and test loss, enhanced training stability, and more easily decodable world features. We also observed reduced distribution drift in certain cases for models trained with the linear probe. The scaling laws for both model size and training time maintained the same slope both with and without the additional of linear probe loss, consistent with findings in the literature on scaling laws for neural networks (Kaplan et al., 2020). However, the networks trained with the linear probes had the loss curve shifted down, such that, for example, a network of a given size trained with the linear probe achieved performance roughly equivalent to a network twice the size trained without the probe. This technique may prove especially advantageous in compute-limited settings or when aiming to achieve the best performance with smaller models.

Our findings support the annealing technique, wherein the linear probe loss term is applied only during the initial stages of training and gradually diminished to zero. The analysis of the loss curve (Figure 6) shows that the slope stabilizes after a certain number of steps, indicating that the linear probes have the most pronounced impact during the early stages of training when the network is actively learning its foundational representations. This suggests that the technique might be particularly beneficial at the beginning of training and less critical in later stages. By focusing the influence of the probes early on, this approach enhances the decodability of world features without imposing a lasting computational burden. Furthermore, transitioning to normal training procedures after the

378 initial phase would remove the additional data requirement, making this technique both efficient and
 379 scalable. Testing this phased application of the technique is a promising direction for future work.
 380

381 Moreover, the concept of leveraging additional sensory inputs aligns with applications in robotic
 382 systems. When a robot is augmented with inexpensive, additional sensors, decoding these inputs
 383 through the network’s latent representations can enhance its ability to interpret and model the en-
 384 vironment. This process facilitates improved learning by enriching the internal state of the model
 385 with diverse sensory information. This technique is thus promising for scalable and cost-effective
 386 advances in robotics, where this addition of sensors can unlock significant performance gains.
 387

388 Overall, this study contributes to our understanding of how to encourage the development of robust
 389 world models in end-to-end trained predictive networks. Further research in this direction may yield
 390 valuable insight for advancing the field of artificial intelligence and developing more capable agents.
 391

392 REFERENCES

393 Christopher M. Bishop. Mixture density networks. Technical Report NCRG/94/004, Neural Com-
 394 puting Research Group, Dept. of Computer Science and Applied Mathematics, Aston University,
 395 Birmingham, B4 7ET, U.K., February 1994. URL <http://www.ncrg aston.ac.uk/>.
 Previously issued as NCRG/94/4288.

396 Alex Graves. Generating sequences with recurrent neural networks, 2014. URL <https://arxiv.org/abs/1308.0850>.

397 David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In
 398 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.),
 399 *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.,
 400 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/2de5d16682c3c35007e4e92982f1a2ba-Paper.pdf.

401 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning
 402 behaviors by latent imagination, 2020. URL <https://arxiv.org/abs/1912.01603>.

403 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
 404 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
 405 models, 2020. URL <https://arxiv.org/abs/2001.08361>.

406
 407
 408
 409
 410
 411
 412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431