# Biologically inspired architectures for sample-efficient deep reinforcement learning

**Pierre H. Richemond**
Imperial College London
phr17@imperial.ac.uk

**Arinbjörn Kolbeinsson**
Imperial College London
ak711@imperial.ac.uk

**Yike Guo**
Imperial College London
y.guo@imperial.ac.uk

## Abstract

Deep reinforcement learning requires a heavy price in terms of sample efficiency and overparameterization in the neural networks used for function approximation. In this work, we use *tensor factorization* in order to learn more compact representation for reinforcement learning policies. We show empirically that in the low-data regime, it is possible to learn online policies with 2 to 10 times less total coefficients, with little to no loss of performance. We also leverage progress in second order optimization, and use the theory of *wavelet scattering* to further reduce the number of learned coefficients, by foregoing learning the topmost convolutional layer filters altogether. We evaluate our results on the Atari suite against recent baseline algorithms that represent the state-of-the-art in data efficiency, and get comparable results with an order of magnitude gain in weight parsimony.

## 1 Introduction & Related Work

The successes of deep reinforcement learning (thereafter 'RL') come at a heavy computational price. It is well known that achieving human-level performance in domains such as Atari [1, 2, 3] requires hundreds of millions of frames of environment interaction. As such, the problem of sample efficiency in reinforcement learning is of critical importance. Several tracks of concurrent research are being investigated, and have reduced by orders of magnitude the number of environment interactions required for good performance beyond the previous benchmark of biologically-inspired episodic control methods [4, 5] to a couple hours of human gameplay time [6, 7].

However, while the data-efficiency of RL methods has seen recent drastic performance, their function approximators still use millions of learned weights, potentially still leaving them heavily overparameterized. Independently motivated by biological facts like the behavioural readiness of newborn animals, several authors [8, 9, 10] have recently looked at doing away with learning so many weights for RL tasks. Smaller networks not only train faster, but may yet offer another avenue for gains in the form of better generalization [11]. Very recent work from [8] studies the effect of inductive bias of neural architectures in reinforcement learning ; they forego training altogether, but transfer networks that only obtain 'better than chance performance on MNIST'. In similar fashion, [10] investigate the effect of random projections in the restricted setting of imitation learning. Finally, [9] manage human-level performance on the Atari suite using a separate dictionary learning procedure for their features, bypassing the usual end-to-end learning paradigm. The perspective of neural architecture search applied to RL appears difficult, if not computationally inextricable.

Concurrently, the study of biologically-inspired models of learning has exhibited two mathematical characterizations that might be critical in explaining how biological learning takes place so efficiently. First, the low-rank properties of learned perceptual manifolds [12, 13] are giving rise to a rich theory borrowing from statistical physics. Second, another well known line of work has identified Gabor filters (and more generally wavelet filter-like structures) in the actual visual cortex of animals [14], and linked those to sparsity-promoting methods and dictionary learning [15, 16, 17]. But these

breakthroughs have not, so far, been reflected as inductive priors in the shape of modifications in deep RL neural networks architectures, which remain fairly fixed on the Atari domain.

Therefore the following questions remain: how parsimonious do function approximators in reinforcement learning need to be, in order to maintain good performance? And can we be at once *sample-efficient* and *weight-efficient* ? In this work, we turn to the mathematical theories of tensor factorization [18], second-order optimization [19, 20] and wavelet scattering [21] to answer this question positively and empirically, in a model-free setting. To the best of our knowledge, this is the first time those fields have been combined together in this context, and that tensor factorization is applied to deep RL.

## 2 Background

### 2.1 Deep Reinforcement Learning

We consider the standard Markov Decision Process framework as in [1]). This setting is characterised by a tuple $\langle S, A, T, R, \gamma \rangle$, where $S$ is a set of states, $A$ a set of actions, $R$ a reward function that is the immediate, intrinsic desirability of a certain state, $T$ a transition dynamics and $\gamma \in [0, 1]$ a discount factor. The purpose of the RL problem is to to find a policy $\pi$, which represents a mapping from states to a probability distribution over actions, that is optimal, i.e., that maximizes the expected cumulative discounted return $\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ at each state $s_t \in S$. In Q-learning, the policy is given implicitly by acting greedily or $\epsilon$-greedily with respect to learned *action-value functions* $q^\pi(s, a)$, that are learned following the Bellman equation. In *deep Q-learning*, $q_\theta$ becomes parameterized by the weights $\theta$ of a neural network and one minimizes the expected Bellman loss :

$$\mathbb{E}\left(R_{t+1} + \gamma_{t+1} \max_{a'} q_\theta\left(S_{t+1}, a'\right) - q_\theta\left(S_t, A_t\right)\right)^2$$

In practice, this is implemented stochastically via uniform sampling of transitions in an experience replay buffer, as is done in the seminal paper [2]. Several algorithmic refinements to that approach exist. First, Double Q-learning [22] proposes to decouple learning between two networks in order to alleviate the Q-value overestimation problem. Second, *dueling* Q-networks [23] explicitly decompose the learning of an action-value function $q_\theta(s, a)$ as the sum of an action-independent state-value, much like what is traditionally done in policy gradient methods [1], implemented via a two-headed neural network architecture. Finally, *prioritized* RL [24] proposes to replace the uniform sampling of transitions in the experience replay buffer with importance sampling, by prioritizing those transitions that present the most Bellman error (those transitions that are deemed the most 'surprising' by the agent). [25] uses extra weights to learn the variance of the exploration noise in a granular fashion, while [26] proposes to learn a full *distribution* of action-values for each action and state. Combined, those methods form the basis of the Rainbow algorithm in [3].

### 2.2 Tensor factorization

Here we introduce notations and concepts from the tensor factorization literature. An intuition is that the two main decompositions below, *CP* and *Tucker* decompositions, can be understood as multilinear algebra analogues of SVD or eigendecomposition.

**CP decomposition.** A tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$, can be decomposed into a sum of $R$ rank-1 tensors, known as the Canonical-Polyadic decomposition, where $R$ is as the rank of the decomposition. The objective is to find the vectors $\mathbf{u}_k^{(1)}, \mathbf{u}_k^{(2)}, \cdots, \mathbf{u}_k^{(N)}$, for $k = [1 \ldots R]$, as well as a vector of weights $\boldsymbol{\lambda} \in \mathbb{R}^R$ such that:

$$\mathcal{X} = \sum_{k=1}^{R} \lambda_k \underbrace{\mathbf{u}_k^{(1)} \circ \mathbf{u}_k^{(2)} \circ \cdots \circ \mathbf{u}_k^{(N)}}_{\text{rank-1 components}}$$

**Tucker decomposition.** A tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$, can be decomposed into a low rank approximation consisting of a core $\mathcal{G} \in \mathbb{R}^{R_1 \times R_2 \times \cdots \times R_N}$ and a set of projection factors $\left(\mathbf{U}^{(0)}, \cdots, \mathbf{U}^{(N-1)}\right)$, with $\mathbf{U}^{(k)} \in \mathbb{R}^{R_k, \hat{I}_k}, k \in (0, \cdots, N-1)$ that, when projected along the corresponding dimension of

the core, reconstruct the full tensor $\mathcal{X}$. The tensor in its decomposed form can then be written:

$$\mathcal{X} = \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times \cdots \times_N \mathbf{U}^{(N)} = \left[ \mathcal{G}; \mathbf{U}^{(1)}, \cdots, \mathbf{U}^{(N)} \right]$$

**Tensor regression layer.** For two tensors $\mathcal{X} \in \mathbb{R}^{K_1 \times \cdots \times K_x \times I_1 \times \cdots \times I_N}$ and $\mathcal{Y} \in \mathbb{R}^{I_1 \times \cdots \times I_N \times L_1 \times \cdots \times L_y}$, we denote by $\langle \mathcal{X}, \mathcal{Y} \rangle_N \in \mathbb{R}^{K_1 \times \cdots \times K_x \times L_1 \times \cdots \times L_y}$ the contraction of $\mathcal{X}$ by $\mathcal{Y}$ along their $N$ last modes; their generalized inner product is

$$\langle \mathcal{X}, \mathcal{Y} \rangle_N = \sum_{i_1=1}^{I_1} \sum_{i_2=1}^{I_2} \cdots \sum_{i_n=1}^{I_N} \mathcal{X}_{\ldots, i_1, i_2, \ldots, i_n} \mathcal{Y}_{i_1, i_2, \ldots, i_n, \ldots}$$

This enables us to define a *tensor regression layer* [27] that is differentiable and learnable end-to-end by gradient descent. Let us denote by $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$ the input activation tensor for a sample and $\mathbf{y} \in \mathbb{R}^{I_N}$ the label vector. A tensor regression layer estimates the regression weight tensor $\mathcal{W} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$ under a low-rank decomposition. In the case of a Tucker decomposition (as per our experiments) with ranks $(R_1, \cdots, R_N)$, we have :

$$\mathbf{y} = \langle \mathcal{X}, \mathcal{W} \rangle_N + \mathbf{b} \qquad \text{with } \mathcal{W} = \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \cdots \times_N \mathbf{U}^{(N)}$$

as $\mathcal{G} \in \mathbb{R}^{R_1 \times \cdots \times R_N}$, $\mathbf{U}^{(k)} \in \mathbb{R}^{I_k \times R_k}$ for each $k$ in $[1 \ldots N]$ and $\mathbf{U}^{(N)} \in \mathbb{R}^{1 \times R_N}$.

[28, 27, 29] learn parsimonious deep learning fully-connected layers thanks to low-rank constraints.

## 2.3 Wavelet scattering

The *wavelet scattering transform* was originally introduced by [21] and [30] as a non-linear extension to the classical wavelet filter bank decomposition [31]. Its principle is as follows. Denoting by $x \circledast y[n]$ the 2-dimensional, circular convolution of two signals $x[n]$ and $y[n]$, let us assume that we have pre-defined two wavelet filter banks available $\left\{ \psi_{\lambda_1}^{(1)}[n] \right\}_{\lambda_1 \in \Lambda_1}$ $\left\{ \psi_{\lambda_2}^{(2)}[n] \right\}_{\lambda_2 \in \Lambda_2}$, with $\lambda_1$ and $\lambda_2$ two frequency indices. These wavelet filters correspond to high frequencies, so we also give ourselves the data of a lowpass filter $\phi_J[n]$. Finally, and by opposition to traditional linear wavelet transforms, we also assume a given nonlinearity $\rho(t)$. Then the scattering transform is given by coefficients of order 0,1, and 2, respectively :

$$S_0 x[n] = x \circledast \phi_J[n]$$

$$S_1 x[n, \lambda_1] = \rho \left( x \circledast \psi_{\lambda_1}^{(1)} \right) \circledast \phi_J[n] \quad \lambda_1 \in \Lambda_1$$

$$S_2 x[n, \lambda_1, \lambda_2] = \rho \left( \rho \left( x \circledast \psi_{\lambda_1}^{(1)} \right) \circledast \psi_{\lambda_2}^{(2)} \right) \circledast \phi_J[n] \quad \lambda_1 \in \Lambda_1, \lambda_2 \in \Lambda_2(\lambda_1)$$

This can effectively be understood and implemented as a two-layer convolutional neural network whose weights are not learned but rather frozen and given by the coefficients of wavelets $\psi$ and $\phi$ (with Gabor filters as a special case [31]). The difference with traditional filter banks comes from the iterated modulus/nonlinear activation function applied at each stage, much like in traditional deep learning convolutional neural networks. In practice, the potential of scattering transforms to accelerate deep learning by providing ready-made convolutional layers weights has been investigated in [32, 33, 34].

## 2.4 Second order optimization with K-FAC

While stochastic gradient descent is usually performed purely from gradient observations derived from auto-differentiation, faster, second order optimization methods first multiply the weights' $\theta$ gradient vector $\nabla_\theta$ by a preconditioning matrix, yielding the weight update $\theta \leftarrow \theta - \eta G^{-1} \nabla_\theta$. In the case of second order methods, the matrix $G^{-1}$ is chosen to act as a tractable iterative approximation to the inverse Hessian or Empirical Fisher Information Matrix [19] of the neural network model in question. Kronecker-factored approximate curvature or K-FAC [20] enforces a Kronecker decomposition of the type $G = A \otimes B$, with $A$ and $B$ being smaller, architecture-dependent matrices. Unlike the above methods, K-FAC *has* been applied as a plug-in in the RL literature and been shown to promote convergence [35].

# 3 Methods & Experimental Results

We do take as a baseline method the data-efficient Rainbow of [6]. However, we change the architecture of the neural network function approximators used, in accordance with the principles described above, combining them to reflect inductive biases promoting fewer learnable parameters:

- We replace the fully-connected, linear layers used in the Rainbow [3] and data-efficient Rainbow [6] by tensor regression layers [27] in order to learn *low-rank* policies (ranks in appendix).
- We use either the K-FAC [20] second order stochastic optimizer, or ADAM [36].
- We combine the two methods with various rank and therefore weight compression ratios and evaluate those on the same subset of Atari games as [6, 7].
- When possible, we replace the first convolutional layer in the approximating neural network with a *scattering* layer for further gains in terms of learnable weights.

For all our Atari experiments, we used OpenAI Gym [37], and a combination of PyTorch [38], TensorLy [39] and Kymatio [40] for auto-differentiation. We evaluated our agents in the low-data regime of 100,000 steps, on half the games, with 3 different random seeds for reproducibility [41]. Our specific hyperparameters are described in appendix. We report our results in tables 1 and 2.

| Game | SimPLe | Rainbow | Denoised | TRL 2.5x | TRL 5x | TRL 10x |
|---|---|---|---|---|---|---|
| alien | 405 | **740** | 684 | 688 | 454 | 566 |
| amidar | 88 | **189** | 154 | 118 | 86 | 84 |
| assault | 369 | 431 | 321 | **543** | 521 | 513 |
| asterix | **1090** | 471 | 500 | 459 | 554 | 363 |
| bank_heist | 8 | 51 | 77 | 59 | **134** | 42 |
| battle_zone | 5184 | 10125 | 9378 | **14466** | 13466 | 5744 |
| boxing | **9** | 0.2 | 1 | -2 | -2 | -5 |
| breakout | **13** | 2 | 3 | 2 | 2 | 4 |
| chopper_command | 1247 | 862 | **1293** | 1255 | 1243 | 1106 |
| crazy_climber | **39828** | 16185 | 9977 | 3928 | 4225 | 2340 |
| demon_attack | 170 | **508** | 450 | 362 | 263 | 175 |
| freeway | 20 | **28** | 28 | 26 | 25 | 24 |
| frostbite | 255 | 867 | **1101** | 659 | 912 | 231 |
| gopher | **771** | 349 | 391 | 278 | 255 | 396 |
| hero | 1295 | **6857** | 3013 | 5351 | 3732 | 3321 |
| jamesbond | 125 | **302** | 295 | 215 | 213 | 218 |
| kangaroo | 323 | **779** | 1002 | 804 | 715 | 400 |
| krull | **4540** | 2852 | 2656 | 2333 | 2275 | 2308 |
| kung_fu_master | **17257** | 14346 | 4037 | 9392 | 4764 | 4031 |
| ms_pacman | 763 | **1204** | 1053 | 818 | 838 | 517 |
| pong | **5** | -19 | -20 | -20 | -19 | -21 |
| private_eye | 58 | 98 | 100 | 51 | 100 | **1128** |
| qbert | 560 | **1153** | 672 | 697 | 581 | 733 |
| road_runner | 5169.4 | **9600** | 5426 | 6965 | 3914 | 1319 |
| seaquest | 371 | 354 | **387** | 345 | 350 | 287 |
| up_n_down | 2153 | 2877 | **5123** | 2197 | 2302 | 2179 |
| **Average (vs. Rainbow)** | | **100%** | **118%** | **96%** | **90%** | **71%** |

Table 1: Mean episode returns as reported in SimPLe [7] and data-efficient Rainbow [6], versus our agents, on 26 Atari games. *'Denoised'* is the NoisyNet ablation of Rainbow; *'TRL'* shows the performance of the data-efficient Rainbow with tensor regression layers substituted for linear ones.

Table 1 shows proof of concept of the online learning of low-rank policies, with a loss of final performance varying in proportion to the compression in the low-rank linear layers, very much like in the deep learning literature [28, 27]. The number of coefficients in the original data-efficient Rainbow

is of the order of magnitude of 1M and varies depending on the environment and its action-space size. The corresponding tensor regression layer ranks are in appendix, and chosen to target 400k, 200k and 100k coefficients respectively. While individual game results tend to decrease monotonously with increasing compression, we observe that they are noisy as per the nature of exploration in RL, and average scores reported correspond with the intuition that performance seems to decrease fast after a certain overparameterization threshold is crossed. To take this noisy character into account, we take care to be conservative and report the average of the final three episodes of the learned policy after 80000, 90000 and 100000 steps, respectively. Also, so as to not muddy the discussion and provide fair baselines, we do report on the NoisyNet [25] ablation of Rainbow ('Denoised' columns), as the NoisyLinear layer doubles up the number of coefficients required and actually performs worse in our experiments. Interestingly, the approximation error in tensor factorization seems to play a role akin to promoting exploration noise.

We then proceed to assess the impact of second-order optimization to our architecture by substituting ADAM optimization for K-FAC, and introducing scattering, in table 2 (only a handful results being available with scattering, due to computational limitations). In spite of our conservative reporting, the efficiency boost from using a second order scheme more than makes up for low-rank approximation error with five times less coefficients than [6], suggesting that learning with a full order of magnitude less coefficients is well within reach of our techniques.

| Game | KFAC+Denoised | KFAC+TRL5x | KFAC+TRL10x | Scattering |
|---|---|---|---|---|
| alien | **996** | 734 | 643 | 441 |
| amidar | **163** | 101 | 98 | 84 |
| assault | **501** | 491 | 496 | 434 |
| asterix | 537 | **549** | 526 | 502 |
| bank_heist | **100** | 73 | 57 | 29 |
| battle_zone | 8622 | **15178** | 6156 | 4311 |
| boxing | **0** | -4 | -1 | -9 |
| breakout | 3 | **3** | 2 | 2 |
| chopper_command | 692 | 611 | **1302** | 441 |
| crazy_climber | **14242** | 12377 | 3546 | 740 |
| demon_attack | 582 | 434 | 318 | **692** |
| freeway | 26 | 26 | 24 | 19 |
| frostbite | **1760** | 718 | 1483 | 654 |
| gopher | **363** | 341 | 265 | 172 |
| hero | 4188 | **6284** | 4206 | 4127 |
| jamesbond | 263 | **327** | 217 | 48 |
| kangaroo | **2085** | 613 | 588 | 391 |
| krull | 2855 | **3441** | 3392 | 772 |
| kung_fu_master | 8481 | **10738** | 7357 | 233 |
| ms_pacman | **1137** | 920 | 867 | 613 |
| pong | -19.3 | **-19** | -19 | -20 |
| private_eye | 56 | **100** | 100 | 0 |
| qbert | **731** | 520 | 538 | 475 |
| road_runner | 4516 | **8493** | 7224 | 1278 |
| seaquest | 349 | 317 | **520** | 213 |
| up_n_down | **2557** | 2291 | 2108 | 993 |
| **Average (vs. Rainbow)** | **114%** | **109%** | **98%** | **56%** |

Table 2: Mean episode returns of our low-rank agents with second-order optimization and scattering.

## 4 Conclusion

We have demonstrated that in the low-data regime, it is possible to leverage biologically plausible characterizations of experience data (namely low-rank properties and wavelet scattering separability) to exhibit architectures that learn policies with many times less weights than current baselines, *in an online fashion*. We do hope that this will lead to even further progress towards sample efficiency and speedy exploration methods. Further work will first focus on thorough evaluation and research of scattering architectures in order to achieve further gains, and second investigate additional, orthogonal biologically-friendly research directions such as promoting sparsity.

# References

[1] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. The MIT Press, second ed., 2018.

[2] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing Atari with Deep Reinforcement Learning," *arXiv e-prints*, Dec. 2013.

[3] M. Hessel, J. Modayil, H. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. Azar, and D. Silver, "Rainbow: Combining Improvements in Deep Reinforcement Learning," *arXiv e-prints*, Oct. 2017.

[4] C. Blundell, B. Uria, A. Pritzel, Y. Li, A. Ruderman, J. Z. Leibo, J. Rae, D. Wierstra, and D. Hassabis, "Model-Free Episodic Control," *arXiv e-prints*, June 2016.

[5] A. Pritzel, B. Uria, S. Srinivasan, A. Puigdomènech, O. Vinyals, D. Hassabis, D. Wierstra, and C. Blundell, "Neural Episodic Control," *arXiv e-prints*, Mar. 2017.

[6] H. van Hasselt, M. Hessel, and J. Aslanides, "When to use parametric models in reinforcement learning?," *arXiv e-prints*, June 2019.

[7] L. Kaiser, M. Babaeizadeh, P. Milos, B. Osinski, R. H. Campbell, K. Czechowski, D. Erhan, C. Finn, P. Kozakowski, S. Levine, A. Mohiuddin, R. Sepassi, G. Tucker, and H. Michalewski, "Model-Based Reinforcement Learning for Atari," *arXiv e-prints*, Mar. 2019.

[8] A. Gaier and D. Ha, "Weight Agnostic Neural Networks," *arXiv e-prints*, June 2019.

[9] G. Cuccu, J. Togelius, and P. Cudre-Mauroux, "Playing Atari with Six Neurons," *arXiv e-prints*, June 2018.

[10] R. Wang, C. Ciliberto, P. Amadori, and Y. Demiris, "Random Expert Distillation: Imitation Learning via Expert Policy Support Estimation," *arXiv e-prints*, May 2019.

[11] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," *arXiv e-prints*, Nov. 2016.

[12] S. Chung, D. D. Lee, and H. Sompolinsky, "Classification and Geometry of General Perceptual Manifolds," *Physical Review X*, vol. 8, p. 031003, July 2018.

[13] S. Chung, D. D. Lee, and H. Sompolinsky, "Linear readout of object manifolds.," *Physical review. E*, vol. 93 6, p. 060301, 2016.

[14] J. P. Jones and L. A. Palmer, "An evaluation of the two-dimensional gabor filter model of simple receptive fields in cat striate cortex.," *Journal of neurophysiology*, vol. 58 6, pp. 1233–58, 1987.

[15] B. A. Olshausen and D. J. Field, "Emergence of simple-cell receptive field properties by learning a sparse code for natural images," *Nature*, vol. 381, pp. 607–609, 1996.

[16] B. A. Olshausen and D. J. Field, "Sparse coding with an overcomplete basis set: A strategy employed by v1?," *Vision Research*, vol. 37, pp. 3311–3325, 1997.

[17] A. Hyvärinen and P. O. Hoyer, "A two-layer sparse coding model learns simple and complex cell receptive fields and topography from natural images," *Vision Research*, vol. 41, pp. 2413–2423, 2001.

[18] A. Cichocki, R. Zdunek, A. H. Phan, and S. Amari, *Nonnegative Matrix and Tensor Factorizations - Applications to Exploratory Multi-way Data Analysis and Blind Source Separation*. Wiley, 2009.

[19] S. ichi Amari, "Natural gradient works efficiently in learning," *Neural Computation*, vol. 10, pp. 251–276, 1998.

[20] J. Martens and R. B. Grosse, "Optimizing neural networks with kronecker-factored approximate curvature," *ArXiv*, vol. abs/1503.05671, 2015.

[21] S. Mallat, "Group Invariant Scattering," *arXiv e-prints*, Jan. 2011.

[22] H. van Hasselt, A. Guez, and D. Silver, "Deep Reinforcement Learning with Double Q-learning," *arXiv e-prints*, Sept. 2015.

[23] Z. Wang, T. Schaul, M. Hessel, H. van Hasselt, M. Lanctot, and N. de Freitas, "Dueling Network Architectures for Deep Reinforcement Learning," *arXiv e-prints*, Nov. 2015.

[24] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized Experience Replay," *arXiv e-prints*, Nov. 2015.

[25] M. Fortunato, M. Gheshlaghi Azar, B. Piot, J. Menick, I. Osband, A. Graves, V. Mnih, R. Munos, D. Hassabis, O. Pietquin, C. Blundell, and S. Legg, "Noisy Networks for Exploration," *arXiv e-prints*, June 2017.

[26] M. G. Bellemare, W. Dabney, and R. Munos, "A Distributional Perspective on Reinforcement Learning," *arXiv e-prints*, July 2017.

[27] J. Kossaifi, Z. C. Lipton, A. Khanna, T. Furlanello, and A. Anandkumar, "Tensor Regression Networks," *arXiv e-prints*, July 2017.

[28] J. Kossaifi, A. Khanna, Z. C. Lipton, T. Furlanello, and A. Anandkumar, "Tensor Contraction Layers for Parsimonious Deep Nets," *arXiv e-prints*, June 2017.

[29] X. Cao and G. Rabusseau, "Tensor Regression Networks with various Low-Rank Tensor Approximations," *arXiv e-prints*, Dec. 2017.

[30] J. Bruna and S. Mallat, "Invariant Scattering Convolution Networks," *arXiv e-prints*, Mar. 2012.

[31] S. Mallat, *A Wavelet Tour of Signal Processing*. Academic Press, 1998.

[32] E. Oyallon, S. Mallat, and L. Sifre, "Generic Deep Networks with Wavelet Scattering," *arXiv e-prints*, Dec. 2013.

[33] E. Oyallon, S. Zagoruyko, G. Huang, N. Komodakis, S. Lacoste-Julien, M. Blaschko, and E. Belilovsky, "Scattering Networks for Hybrid Representation Learning," *arXiv e-prints*, Sept. 2018.

[34] T. Angles and S. Mallat, "Generative networks as inverse problems with Scattering transforms," *arXiv e-prints*, May 2018.

[35] Y. Wu, E. Mansimov, S. Liao, R. Grosse, and J. Ba, "Scalable trust-region method for deep reinforcement learning using Kronecker-factored approximation," *arXiv e-prints*, Aug. 2017.

[36] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv e-prints*, Dec. 2014.

[37] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "OpenAI Gym," *arXiv e-prints*, June 2016.

[38] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *OpenReview*, 2017.

[39] J. Kossaifi, Y. Panagakis, A. Anandkumar, and M. Pantic, "TensorLy: Tensor Learning in Python," *arXiv e-prints*, Oct. 2016.

[40] M. Andreux, T. Angles, G. Exarchakis, R. Leonarduzzi, G. Rochette, L. Thiry, J. Zarka, S. Mallat, J. andén, E. Belilovsky, J. Bruna, V. Lostanlen, M. J. Hirn, E. Oyallon, S. Zhang, C. Cella, and M. Eickenberg, "Kymatio: Scattering Transforms in Python," *arXiv e-prints*, Dec. 2018.

[41] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, "Deep Reinforcement Learning that Matters," *arXiv e-prints*, Sept. 2017.

[42] M. Eickenberg, G. Exarchakis, M. J. Hirn, S. Mallat, and L. Thiry, "Solid harmonic wavelet scattering for predictions of molecule properties," *The Journal of chemical physics*, vol. 148 24, p. 241732, 2018.

# Appendix

## Hyperparameters and Reproducibility

Our codebase is available on request. Hyperparameters are as follows. First, our specific architecture-modified hyperparameters:

| Specific architecture hyperparameters | Value |
|---|---|
| Scattering maximum log-scale $J$ | 3 |
| Scattering volume width $M$ | 1 |
| Scattering tensor input shape | (1,4,84,84) |
| Scattering tensor output shape | (1,16,11,11) |
| Scattering type | Harmonic 3D, see [40, 42] |
| Hidden linear layer rank constraint, 2.5x compression | 128 |
| Final linear layer rank constraint, 2.5x compression | 48 |
| Hidden linear layer rank constraint, 5x compression | 32 |
| Final linear layer rank constraint, 5x compression | 48 |
| Hidden linear layer rank constraint, 10x compression | 16 |
| Final linear layer rank constraint, 10x compression | 10 |
| KFAC Tikhonov regularization parameter | 0.1 |
| KFAC Update frequency for inverses | 100 |

Table 3: Our additional, architecture-specific hyperparameters.

Furthermore, we mirror the Data-Efficient Rainbow [6] baseline:

| Data-efficient Rainbow hyperparameters | Value |
|---|---|
| Grey-scaling | True |
| Observation down-sampling | (84, 84) |
| Frames stacked | 4 |
| Action repetitions | 4 |
| Reward clipping | [-1, 1] |
| Terminal on loss of life | True |
| Max frames per episode | 108K |
| Update | Distributional Double Q |
| Target network update period[*] | every 2000 updates |
| Support of Q-distribution | 51 bins |
| Discount factor | 0.99 |
| Minibatch size | 32 |
| Optimizer | Adam |
| Optimizer: first moment decay | 0.9 |
| Optimizer: second moment decay | 0.999 |
| Optimizer: $\epsilon$ | 0.00015 |
| Max gradient norm | 10 |
| Priority exponent | 0.5 |
| Priority correction[**] | $0.4 \rightarrow 1$ |
| Hardware | NVidia 1080Ti GPU |
| Noisy nets parameter | 0.1 |
| Training frames | 400,000 |
| Min replay size for sampling | 1600 |
| Memory size | unbounded |
| Replay period every | 1 steps |
| Multi-step return length | 20 |
| Q network: channels | 32, 64 |
| Q network: filter size | 5 x 5, 5 x 5 |
| Q network: stride | 5, 5 |
| Q network: hidden units | 256 |
| Optimizer: learning rate | 0.0001 |

Table 4: Data-efficient Rainbow agent hyperparameters, as per [6].

**Standard deviations for score runs**

| Game | Denoised | TRL 2.5x | TRL 10x |
|---|---|---|---|
| alien | $684 \pm 7$ | $688 \pm 123$ | $566 \pm 38$ |
| amidar | $154 \pm 21$ | $118 \pm 12$ | $84 \pm 15$ |
| assault | $321 \pm 224$ | $543 \pm 94$ | $513 \pm 64$ |
| asterix | $500 \pm 124$ | $459 \pm 91$ | $363 \pm 66$ |
| bank_heist | $77 \pm 23$ | $59 \pm 22$ | $42 \pm 2$ |
| battle_zone | $9378 \pm 2042$ | $14466 \pm 2845$ | $5744 \pm 575$ |
| boxing | $1 \pm 2$ | $-2 \pm 1$ | $-5 \pm 1$ |
| breakout | $3 \pm 1.5$ | $2 \pm 1$ | $4 \pm 0.3$ |
| chopper_command | $1293 \pm 445$ | $1255 \pm 215$ | $1106 \pm 124$ |
| crazy_climber | $9977 \pm 3744$ | $3928 \pm 221$ | $2340 \pm 595$ |
| demon_attack | $450 \pm 49$ | $362 \pm 147$ | $175 \pm 7$ |
| freeway | $28 \pm 0.6$ | $26 \pm 0$ | $24 \pm 0.5$ |
| frostbite | $1101 \pm 355$ | $659 \pm 523$ | $231 \pm 1$ |
| gopher | $391 \pm 46$ | $278 \pm 39$ | $396 \pm 24$ |
| hero | $3013 \pm 90$ | $5351 \pm 1948$ | $3321 \pm 598$ |
| jamesbond | $295 \pm 57$ | $215 \pm 42$ | $218 \pm 22$ |
| kangaroo | $1002 \pm 587$ | $804 \pm 289$ | $400 \pm 278$ |
| krull | $2656 \pm 180$ | $2333 \pm 309$ | $2308 \pm 268$ |
| kung_fu_master | $4037 \pm 2962$ | $9392 \pm 6289$ | $4031 \pm 3068$ |
| ms_pacman | $1053 \pm 193$ | $818 \pm 94$ | $517 \pm 38$ |
| pong | $-20 \pm 0.4$ | $-20 \pm 0$ | $-21 \pm 0.1$ |
| private_eye | $100 \pm 0$ | $51 \pm 59$ | $1128 \pm 1067$ |
| qbert | $672 \pm 144$ | $697 \pm 78$ | $733 \pm 291$ |
| road_runner | $5426 \pm 2830$ | $6965 \pm 6569$ | $1319 \pm 216$ |
| seaquest | $387 \pm 24$ | $345 \pm 40$ | $287 \pm 87$ |
| up_n_down | $5123 \pm 3146$ | $2197 \pm 231$ | $2179 \pm 178$ |

Table 5: Standard deviations across seeds for runs presented Table 1.

| Game | KFAC+Denoised | KFAC+TRL10x | Scattering |
|---|---|---|---|
| alien | $996 \pm 180$ | $643 \pm 51$ | $441 \pm 90$ |
| amidar | $163 \pm 15$ | $98 \pm 26$ | $84 \pm 11$ |
| assault | $501 \pm 85$ | $496 \pm 129$ | $434 \pm 304$ |
| asterix | $537 \pm 96$ | $526 \pm 64$ | $502 \pm 91$ |
| bank_heist | $100 \pm 14$ | $57 \pm 36$ | $29 \pm 13$ |
| battle_zone | $8622 \pm 5358$ | $6156 \pm 1951$ | $4311 \pm 1517$ |
| boxing | $0 \pm 2$ | $-1 \pm 3$ | $-9 \pm 12$ |
| breakout | $3 \pm 1$ | $2 \pm 2$ | $2 \pm 0$ |
| chopper_command | $692 \pm 81$ | $1302 \pm 328$ | $441 \pm 80$ |
| crazy_climber | $14242 \pm 2936$ | $3546 \pm 1231$ | $740 \pm 291$ |
| demon_attack | $582 \pm 130$ | $318 \pm 168$ | $692 \pm 232$ |
| freeway | $26 \pm 0$ | $24 \pm 0$ | $19 \pm 1$ |
| frostbite | $1760 \pm 448$ | $1483 \pm 466$ | $654 \pm 709$ |
| gopher | $363 \pm 4$ | $265 \pm 67$ | $172 \pm 3$ |
| hero | $4188 \pm 1635$ | $4206 \pm 1862$ | $4127 \pm 1074$ |
| jamesbond | $263 \pm 22$ | $217 \pm 68$ | $48 \pm 10$ |
| kangaroo | $2085 \pm 2055$ | $588 \pm 5$ | $391 \pm 52$ |
| krull | $2855 \pm 156$ | $3392 \pm 2205$ | $772 \pm 560$ |
| kung_fu_master | $8481 \pm 8270$ | $7357 \pm 9200$ | $233 \pm 205$ |
| ms_pacman | $1137 \pm 180$ | $867 \pm 128$ | $613 \pm 159$ |
| pong | $-19 \pm 0.6$ | $-19 \pm 1$ | $-20 \pm 0$ |
| private_eye | $56 \pm 42$ | $100 \pm 0$ | $0 \pm 0$ |
| qbert | $731 \pm 256$ | $538 \pm 114$ | $475 \pm 161$ |
| road_runner | $4516 \pm 2869$ | $7224 \pm 4598$ | $1278 \pm 463$ |
| seaquest | $349 \pm 63$ | $520 \pm 97$ | $213 \pm 96$ |
| up_n_down | $2557 \pm 641$ | $2108 \pm 298$ | $993 \pm 244$ |

Table 6: Standard deviations across seeds for runs presented Table 2.

## Further results and learning curves

**Learning curves** As a simpler version of the experiments in the main text body, we show basic proof of concept on the simple *Pong* Atari game. Our experimental setup consists in using our own implementation of *prioritized double DQN* as a baseline, which combines algorithmic advances from [24] and [22]. We replaced the densely connected layer of the original DQN architecture with a tensor regression layer implementing Tucker decomposition for different Tucker ranks, yielding different network compression factors. (These curves average three different random seeds).

**Qualitative behaviour.** First results, both in terms of learning performance and compression factor, can be seen in figure 1. The two main findings of this experiment are that first, and overall, the final performance of the agent remains unaffected by the tensor factorization, even with high compression rates nearing 10 times. Second, this obviously comes at the expense of stability during training - in tough compression regimes, learning curves are slightly delayed, and their plateauing phases contain occasional noisy drawdowns illustrating the increased difficulty of learning, as seen in figure 2. The extra pathwise noise, however, can be seen as promoting exploration.

Figure 1: Prioritized tensorized DQN on Atari Pong. Original learning curve versus several learning curves for five different Tucker ranks factorizations and therefore parameter compression rates (3 different random seeds each, with a 30 episodes moving average for legibility). Best viewed in colour.
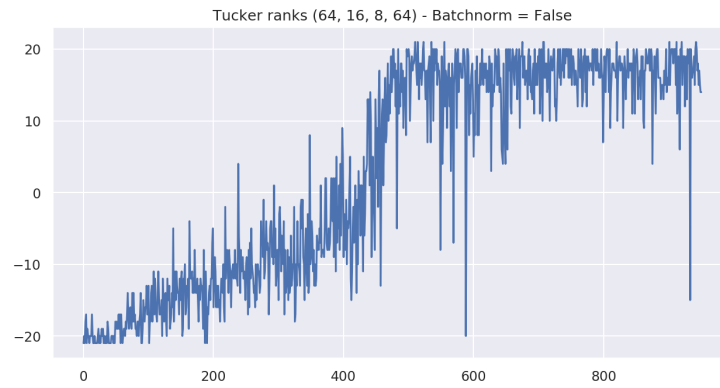


Figure 2: Focus on a typical single run of the tensorized DQN learning. The overall shape of the typical learning curve is preserved, but drawdowns in the plateauing phase do appear.