



**Figure 3 | Comparison of the DQN agent with the best reinforcement learning methods<sup>15</sup> in the literature.** The performance of DQN is normalized with respect to a professional human games tester (that is, 100% level) and random play (that is, 0% level). Note that the normalized performance of DQN, expressed as a percentage, is calculated as:  $100 \times (\text{DQN score} - \text{random play score}) / (\text{human score} - \text{random play score})$ . It can be seen that DQN

outperforms competing methods (also see Extended Data Table 2) in almost all the games, and performs at a level that is broadly comparable with or superior to a professional human games tester (that is, operationalized as a level of 75% or above) in the majority of games. Audio output was disabled for both human players and agents. Error bars indicate s.d. across the 30 evaluation episodes, starting with different initial conditions.

see Fig. 3, Supplementary Discussion and Extended Data Table 2). In additional simulations (see Supplementary Discussion and Extended Data Tables 3 and 4), we demonstrate the importance of the individual core components of the DQN agent—the replay memory, separate target Q-network and deep convolutional network architecture—by disabling them and demonstrating the detrimental effects on performance.

We next examined the representations learned by DQN that underpinned the successful performance of the agent in the context of the game Space Invaders (see Supplementary Video 1 for a demonstration of the performance of DQN), by using a technique developed for the visualization of high-dimensional data called 't-SNE'<sup>225</sup> (Fig. 4). As expected, the t-SNE algorithm tends to map the DQN representation of perceptually similar states to nearby points. Interestingly, we also found instances in which the t-SNE algorithm generated similar embeddings for DQN representations of states that are close in terms of expected reward but

perceptually dissimilar (Fig. 4, bottom right, top left and middle), consistent with the notion that the network is able to learn representations that support adaptive behaviour from high-dimensional sensory inputs. Furthermore, we also show that the representations learned by DQN are able to generalize to data generated from policies other than its own—in simulations where we presented as input to the network game states experienced during human and agent play, recorded the representations of the last hidden layer, and visualized the embeddings generated by the t-SNE algorithm (Extended Data Fig. 1 and Supplementary Discussion). Extended Data Fig. 2 provides an additional illustration of how the representations learned by DQN allow it to accurately predict state and action values.

It is worth noting that the games in which DQN excels are extremely varied in their nature, from side-scrolling shooters (River Raid) to boxing games (Boxing) and three-dimensional car-racing games (Enduro).

**Extended Data Table 2 | Comparison of games scores obtained by DQN agents with methods from the literature<sup>12,15</sup> and a professional human games tester**

Game	Random Play	Best Linear Learner	Contingency (SARSA)	Human	DQN ( $\pm$ std)	Normalized DQN (% Human)
Alien	227.8	939.2	103.2	6875	3069 ( $\pm$ 1093)	42.7%
Amidar	5.8	103.4	183.6	1676	739.5 ( $\pm$ 3024)	43.9%
Assault	222.4	628	537	1496	3359 ( $\pm$ 775)	246.2%
Asterix	210	987.3	1332	8503	6012 ( $\pm$ 1744)	70.0%
Asteroids	719.1	907.3	89	13157	1629 ( $\pm$ 542)	7.3%
Atlantis	12850	62687	852.9	29028	85641 ( $\pm$ 17600)	449.9%
Bank Heist	14.2	190.8	67.4	734.4	429.7 ( $\pm$ 650)	57.7%
Battle Zone	2360	15820	16.2	37800	26300 ( $\pm$ 7725)	67.6%
Beam Rider	363.9	929.4	1743	5775	6846 ( $\pm$ 1619)	119.8%
Bowling	23.1	43.9	36.4	154.8	42.4 ( $\pm$ 88)	14.7%
Boxing	0.1	44	9.8	4.3	71.8 ( $\pm$ 8.4)	1707.9%
Breakout	1.7	5.2	6.1	31.8	401.2 ( $\pm$ 26.9)	1327.2%
Centipede	2091	8803	4647	11963	8309 ( $\pm$ 5237)	63.0%
Chopper Command	811	1582	16.9	9882	6687 ( $\pm$ 2916)	64.8%
Crazy Climber	10781	23411	149.8	35411	114103 ( $\pm$ 22797)	419.5%
Demon Attack	152.1	520.5	0	3401	9711 ( $\pm$ 2406)	294.2%
Double Dunk	-18.6	-13.1	-16	-15.5	-18.1 ( $\pm$ 2.6)	17.1%
Enduro	0	129.1	159.4	309.6	301.8 ( $\pm$ 24.6)	97.5%
Fishing Derby	-91.7	-89.5	-85.1	5.5	-0.8 ( $\pm$ 19.0)	93.5%
Freeway	0	19.1	19.7	29.6	30.3 ( $\pm$ 0.7)	102.4%
Frostbite	65.2	216.9	180.9	4335	328.3 ( $\pm$ 250.5)	6.2%
Gopher	257.6	1288	2368	2321	8520 ( $\pm$ 3279)	400.4%
Gravitar	173	387.7	429	2672	306.7 ( $\pm$ 223.9)	5.3%
H.E.R.O.	1027	6459	7295	25763	19950 ( $\pm$ 158)	76.5%
Ice Hockey	-11.2	-9.5	-3.2	0.9	-1.6 ( $\pm$ 2.5)	79.3%
James Bond	29	202.8	354.1	406.7	576.7 ( $\pm$ 175.5)	145.0%
Kangaroo	52	1622	8.8	3035	6740 ( $\pm$ 2959)	224.2%
Krull	1598	3372	3341	2395	3805 ( $\pm$ 1033)	277.0%
Kung-Fu Master	258.5	19544	29151	22736	23270 ( $\pm$ 5955)	102.4%
Montezuma's Revenge	0	10.7	259	4367	0 ( $\pm$ 0)	0.0%
Ms. Pacman	307.3	1692	1227	15693	2311 ( $\pm$ 525)	13.0%
Name This Game	2292	2500	2247	4076	7257 ( $\pm$ 547)	278.3%
Pong	-20.7	-19	-17.4	9.3	18.9 ( $\pm$ 1.3)	132.0%
Private Eye	24.9	684.3	86	69571	1788 ( $\pm$ 5473)	2.5%
Q*Bert	163.9	613.5	960.3	13455	10596 ( $\pm$ 3294)	78.5%
River Raid	1339	1904	2650	13513	8316 ( $\pm$ 1049)	57.3%
Road Runner	11.5	67.7	89.1	7845	18257 ( $\pm$ 4268)	232.9%
Robotank	2.2	28.7	12.4	11.9	51.6 ( $\pm$ 4.7)	509.0%
Seaquest	68.4	664.8	675.5	20182	5286 ( $\pm$ 1310)	25.9%
Space Invaders	148	250.1	267.9	1652	1976 ( $\pm$ 893)	121.5%
Star Gunner	664	1070	9.4	10250	57997 ( $\pm$ 3152)	598.1%
Tennis	-23.8	-0.1	0	-8.9	-2.5 ( $\pm$ 1.9)	143.2%
Time Pilot	3568	3741	24.9	5925	5947 ( $\pm$ 1600)	100.9%
Tutankham	11.4	114.3	98.2	167.6	186.7 ( $\pm$ 41.9)	112.2%
Up and Down	533.4	3533	2449	9082	8456 ( $\pm$ 3162)	92.7%
Venture	0	66	0.6	1188	380.0 ( $\pm$ 238.6)	32.0%
Video Pinball	16257	16871	19761	17298	42684 ( $\pm$ 16287)	2539.4%
Wizard of Wor	563.5	1981	36.9	4757	3393 ( $\pm$ 2019)	67.5%
Zaxxon	32.5	3365	21.4	9173	4977 ( $\pm$ 1235)	54.1%

Best Linear Learner is the best result obtained by a linear function approximator on different types of hand designed features<sup>12</sup>. Contingency (SARSA) agent figures are the results obtained in ref. 15. Note the figures in the last column indicate the performance of DQN relative to the human games tester, expressed as a percentage, that is,  $100 \times (\text{DQN score} - \text{random play score}) / (\text{human score} - \text{random play score})$ .