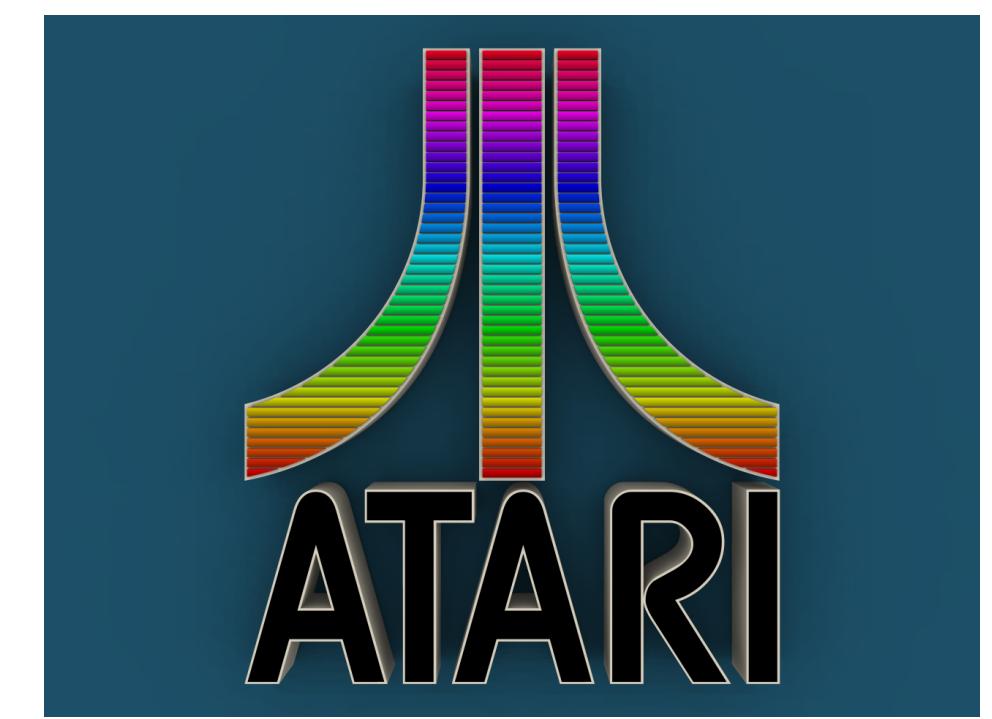


Classical Planning with Simulators: Results on the Atari Video Games



Nir Lipovetzky
The University of Melbourne
Melbourne, Australia

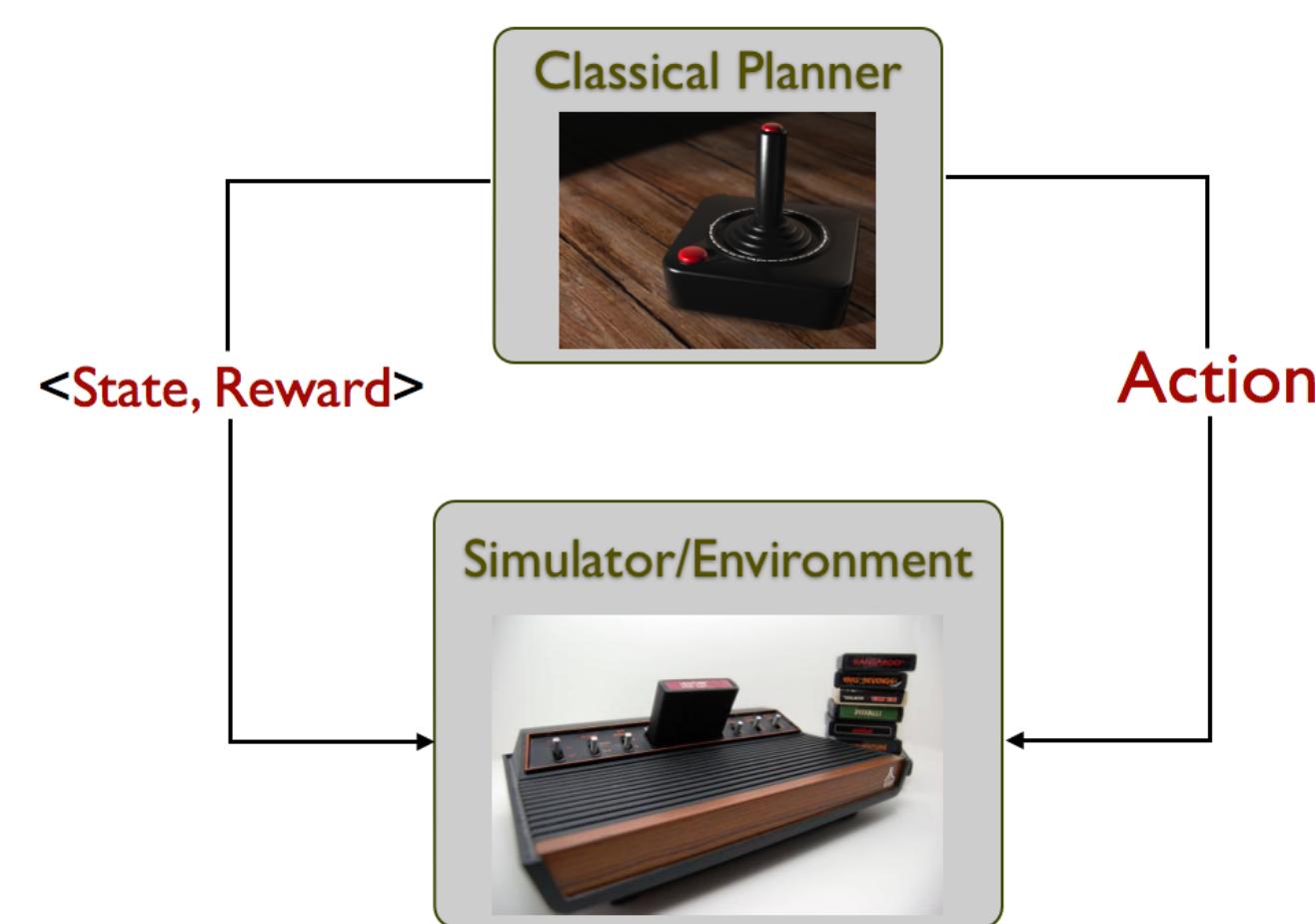
Miquel Ramirez
Australian National University
Canberra, Australia

Hector Geffner
Universitat Pompeu Fabra
Barcelona, Spain

Arcade Learning Environment (U. of Alberta)

Specs of Environment/Simulator:

- RAM: 128 Bytes
- Screen: 160 x 210 pixels with (7 bit) 128-colour palette
- Joystick: 18 actions available
- 60 frames per second. Agents choose action every 5 frames (0.08 sec) of gameplay



ALE and Classical Planning

→ Classical Planning is about achieving goals in deterministic systems with a fully known initial state

→ Problems represented in compact, logical form with PDDL

Planning setting in ALE is deterministic and initial state fully known

Yet classical planners can't be used!

- no PDDL encoding, i.e. No action model available
- no goals but rewards

→ Bellemare et al. consider Breadth-first Search (BrFS) and MCTS

→ Still, some recent "classical" planning algorithms can be applied almost off-the-shelf (Lipovetzky and Geffner 2012)

A Simple Pruned Breadth-First Search Algorithm

Definition (Novelty): smallest new tuple of atoms generated by a state, that no other state in the search made true before. $\text{novelty} : s \mapsto \mathbb{N}$

- $IW(i)$ = breadth-first search that prunes newly generated states whose $\text{novelty}(s) > i$.
- IW is a sequence of calls $IW(i)$ for $i = 0, 1, 2, \dots$ over problem P until problem solved or i exceeds number of vars in problem

Key theoretical property: IW solves P in time $O(n^i)$ if $\text{width}(P) = i$; $IW(i)$ then solves P optimally (Notion of width defined by Lipovetzky and Geffner 2012)

Experiments in a nutshell: IW , while simple and blind, great algorithm when goals restricted to single atoms. For multiple atomic goals, simple serialized version IW suffices, achieving atomic goals one at a time

IW in the Atari Games

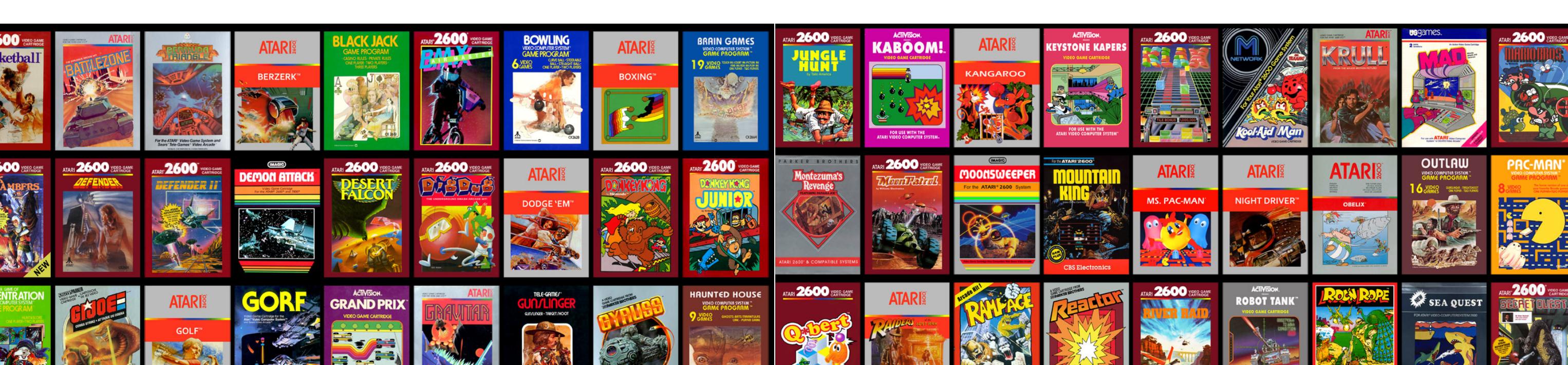
→ $IW(1)$ used with the 128 state variables (RAM 8-bit word) of 256 values each

→ $IW(1)$ generates then up to $128 \times 256 \times 18$ (i.e., 589,824) states

- Children in $IW(1)$ generated in random order
- Discount factor used $\gamma = 0.995$

→ Action leading to most rewarding $IW(1)$ -path is executed

IW is a purely exploration algorithm that does not take into account the accumulated reward for selecting the states to consider



Balancing Exploration and Exploitation

As a simple variant that combines exploration and exploitation, we evaluated a best-first search algorithm with two queues (2BFS)

- Queue 1: ordered first by novelty measure
- Queue 2: ordered by accumulated reward

→ 2BFS alternates between the two queues (similar technique used in LAMA planner)

→ When node expanded, its children are placed on both queues, with the exception of nodes with no accumulated reward, placed only in queue 1

Experimental Results

Same setting from Bellemare et al:

- Games are played for 5 minutes maximum (18,000 frames)
- 2BFS and IW have a maximum lookahead budget of 150,000 simulated frames
- UCT has same budget by running 500 rollouts of depth 300
- Score is averaged among 5 runs per game

Game	$IW(1)$		2BFS		BrFS		UCT	
	Score	Time	Score	Time	Score	Time	Score	Time
ALIEN	25634	81	12525	81	784	7785		
AMIDAR	1377	28	1090	37	5	180		
ASSAULT	953	18	827	25	414	1512		
ASTERIX	153400	24	77200	27	2136	29700		
ASTEROIDS	51338	66	22168	65	3127	4661		
ATLANTIS	159420	13	154180	71	30460	193858		
BANK HEIST	717	39	362	64	22	498		
BATTLE ZONE	11600	86	330800	87	6313	70333		
BEAM RIDER	9108	23	9298	29	694	6625		
BERZERK	2096	58	802	73	195	554		
BOWLING	69	10	50	60	26	25		
BOXING	108	15	100	22	100	100		
BREAKOUT	384	4	772	39	1	364		
CARNIVAL	6372	16	5516	53	950	5132		
CENTIPEDE	99207	39	94236	67	125123	110422		
CHOPPER COMMAND	10980	76	27220	73	1827	34019		
CRAZY CLIMBER	36160	4	36940	58	37110	98172		
DEMON ATTACK	20116	33	16025	41	443	28159		
DOUBLE DUNK	14	41	21	41	19	24		
ELEVATOR ACTION	13480	26	10820	27	730	18100		
ENDURO	500	66	359	38	1	286		
FISHING DERBY	30	39	6	62	-92	38		
FREEWAY	31	32	23	61	0	0		
FROSTBITE	902	12	2672	38	137	271		
GOPHER	18256	19	15808	53	1019	20560		
GRAVITAR	3920	62	5980	62	395	2850		

Game	$IW(1)$		2BFS		BrFS		UCT	
	Score	Time	Score	Time	Score	Time	Score	Time
HERO	12985	37	11524	69	1324	12860		
ICE HOCKEY	55	89	49	89	-	-	39	
JAMES BOND	23070	0	10080	30	25	330		
JOURNEY ESCAPE	40080	38	40600	67	1327	7683		
KANGAROO	8760	8	5320	31	90	1990		
KRULE	6030	28	4884	42	3089	5037		
KUNG FU MASTER	63780	21	42180	43	12127	48855		
MONTEZUMA REVENGE	0	14	540	39	0	0	0	0
MS PACMAN	21695	21	18927	23	1709	22336		
NAME THIS GAME	9354	14	8304	25	5699	15416		
PONG	21	17	21	35	-21	21		
POOYAN	11225	8	10760	16	910	17763		
PRIVATE EYE	-99	18	2544	44	58	100		
Q*BERT	3705	11	11680	35	133	17343		
RIVERRAID	5694	18	5062	37	2179	4449		
ROAD RUNNER	94940	25	68500	41	245	38725		
ROBOT TANK	68	34	52	34	2	50		
SEAQUEST	14272	25	6138	33	288	5132		
SPACE INVADERS	2877	21	3974	34	112	2718		
STAR GUNNER	1540	19	4660	18	1345	1207		
TENNIS	24	21	24	36	-24	3		
TIME PILOT	35000	9	36180	29	4064	63855		
TUTANKHAM	172	15	204	34	64	226		
UP AND DOWN	110036	12	54820	14	746	74474		
VENTURE	1200	22	980	35	0	0		
VIDEO PINBALL	388712	43	62075	43	55567	254748		
WIZARD OF WOR	121060	25	81500	27	3309	105500		
ZAXXON	29240	34	15680	31	0	22610		

	$IW(1)$	2BFS	BrFS	UCT
# Times Best (54 games)	26	13	1	19
# Times Better than IW	—	16	1	19
# Times Better than 2BFS	34	—	1	25
# Times Better than UCT	31	26	1	—

Search Tree Depths

- BrFS search tree results in a lookahead of 0.3 seconds
- $IW(1)$ and 2BFS result in lookahead of up to 6–22 seconds

