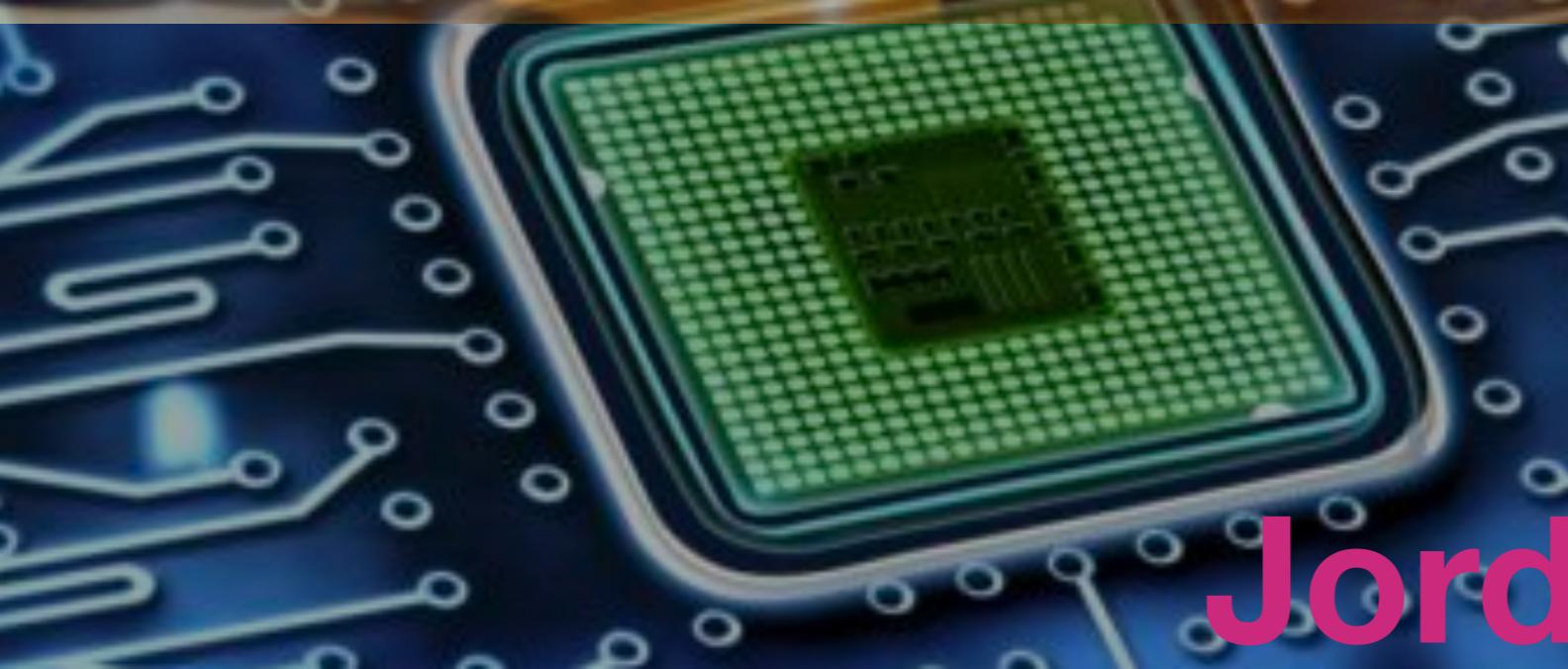


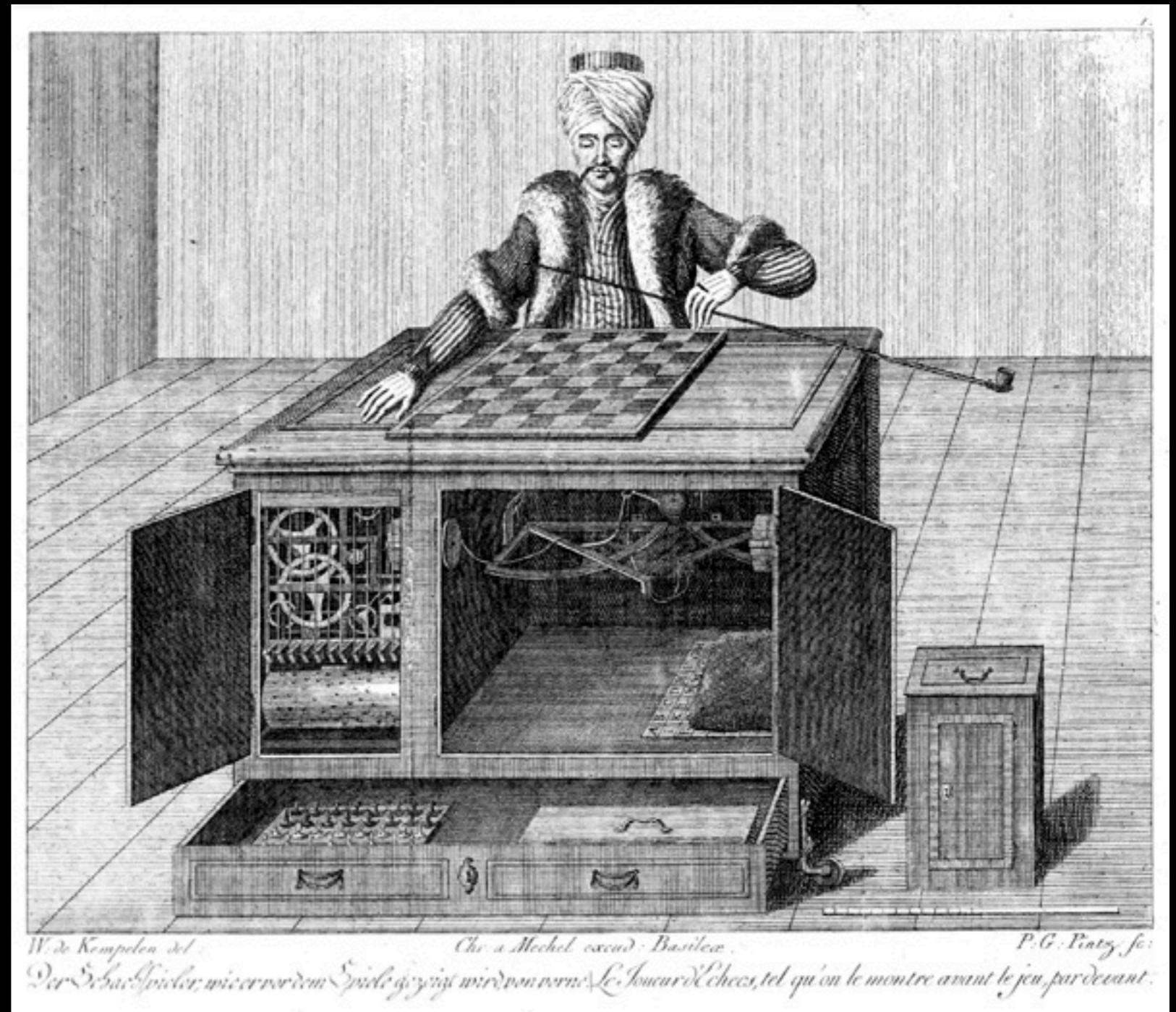
Computers Playing Games



Jordan Hoffmann

1769

- Wolfgang von Kempelen
 - Notable opponents:
 - Napoleon
 - Benjamin Franklin



1868

- Charles Hooper Ajeeb
 - Notable opponents:
 - Harry Houdini
 - Theodore Roosevelt



1878

- Mephisto
 - Controlled remotely



1948

- Turochamp – Alan Turing and David Champernowne
- Written before computers could run it.
- Played Kasparov in 2012 – 2 ply.

1949

- Claude Shannon “Programming a computer for playing chess”

These and similar principles are only generalizations from empirical evidence of numerous games, and only have a kind of statistical validity. Probably any chess principle can be contradicted by particular counter examples. However, from these principles one can construct a crude evaluation function. The following is an example: -

$$f(P) = 200(K-K') + 9(Q-Q') + 5(R-R') + 3(B-B'+N-N') + (P-P') - \\ 0.5(D-D'+S-S'+I-I') + \\ 0.1(M-M') + \dots$$

in which: -

(1)K,Q,R,B,B,P are the number of White kings, queens, rooks, bishops, knights and pawns on the board.

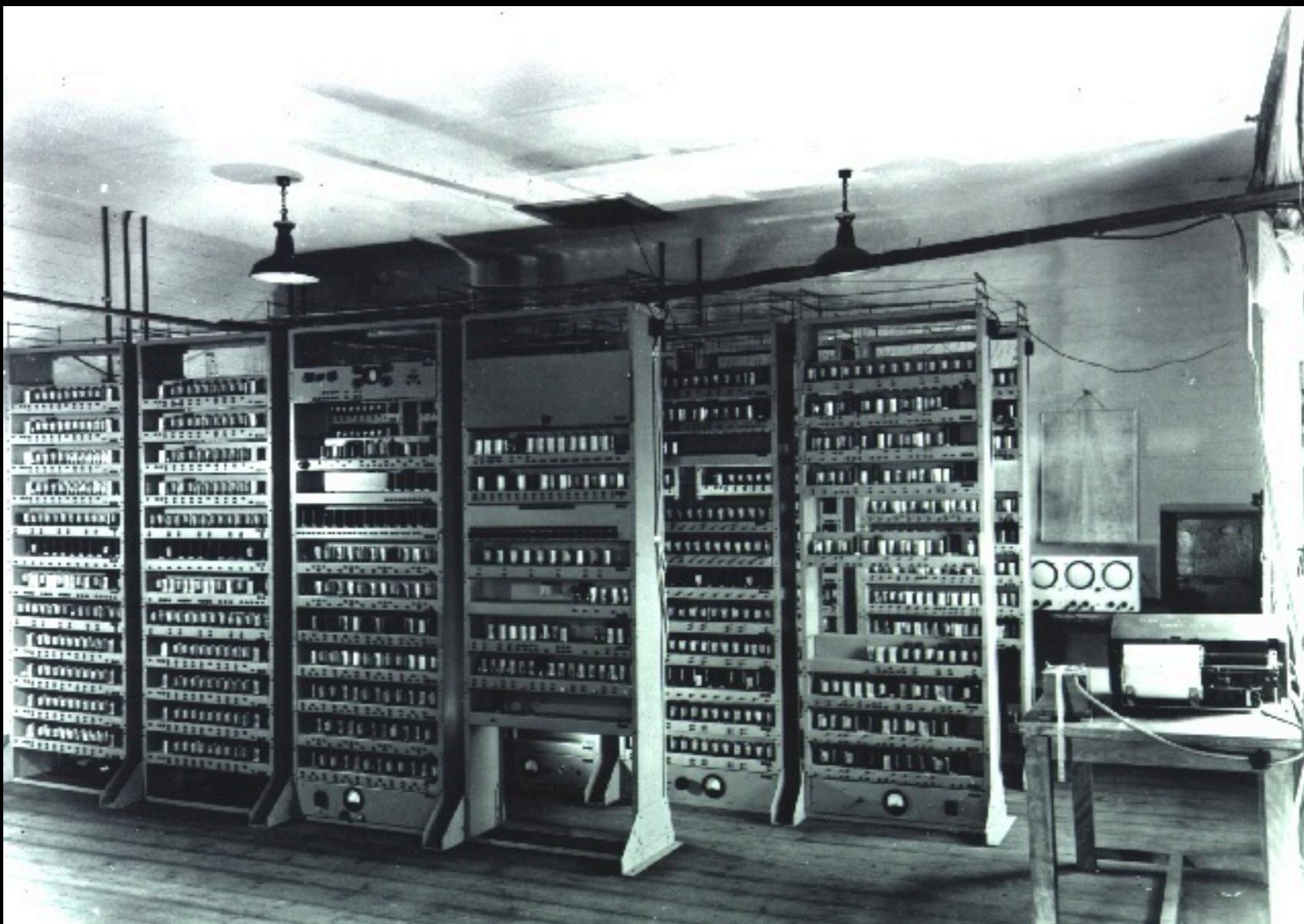
(2)D,S,I are doubled, backward and isolated White pawns.

(3)M= White mobility (measured, say, as the number of legal moves available to White).

Primed letters are the similar quantities for Black.

1952

- OXO PhD thesis at Cambridge, A S Douglas. First computer game.

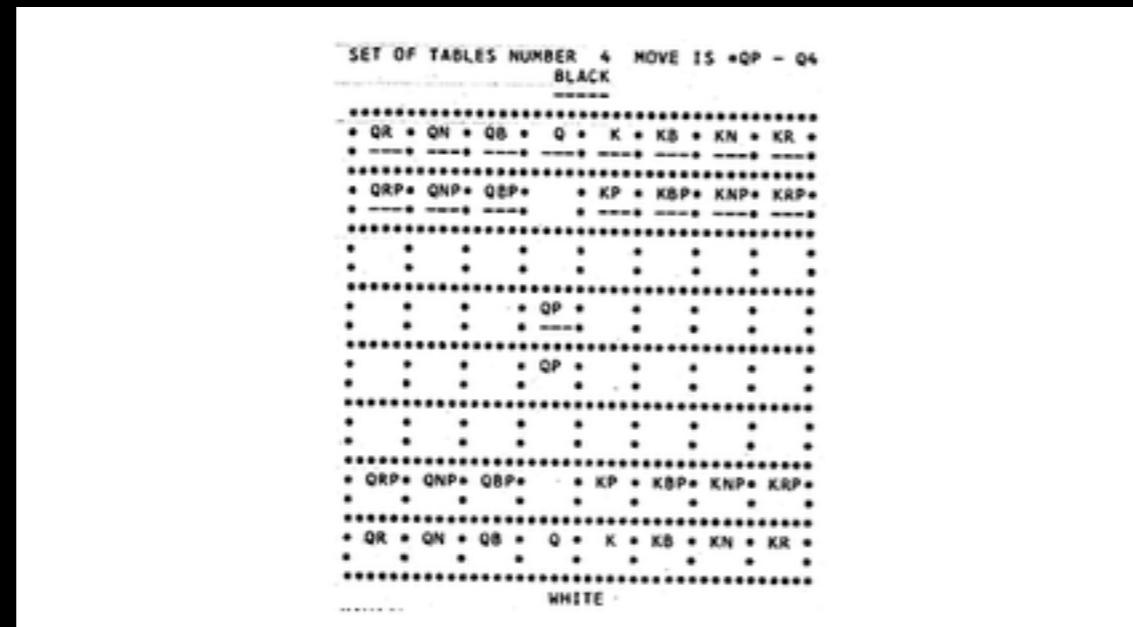


1956

- John McCarthy invents alpha-beta search.

1962

- Kotok-McCarthy is first program to play a recognizable game of chess. *Kotok did this for his MIT undergraduate thesis.*



1986

- Maven— Scrabble, 1998 beat world champion

Find all possible plays

Rank all

Randomly draw tiles for opponent and play two turns

Rack evaluation

1986

- Maven— Scrabble, 1998 beat world champion

Find all possible plays

Rank all

Randomly draw tiles for opponent and play two turns

Rack evaluation

When within one draw of 7 tiles, try to yield good “endgame”

1986

- Maven— Scrabble, 1998 beat world champion

Find all possible plays

Rank all

Randomly draw tiles for opponent and play two turns

Rack evaluation

When within one draw of 7 tiles, try to yield good “endgame”

In the endgame, can use optimal play.

Note: does not use Alpha-Beta (will discuss in a moment)

1986

- Maven— Scrabble, 1998 beat world champion

M	3	O	1	U	1	T	1	H	4		R	1	R	1	T	1			2L		3W					
R	1	E	1							3L					Q	10				2W						
T	1			2W							2L		2L	U	1			G	2							
H	4	U	1	R	1	T	1					2L		R	1		2W	R	1		2L					
N	1	E	1	O	1	N	1						I	1	S	1		E	1	L	1					
	3L		D	2	O	1	Z	10	Y	4			3L	P	3		A	1	X	8	E	1				
	E	1					E	1			2L	J	8	R	1	W	4	S	1		I	1				
I	1	A	1	M	3	B	3		C	3	R	1	V	4	Y	4	N	1	2L	E	1	3W				
W	4	E	1					R	1			2L		K	5				2L							
	3L	N	1		F	4	3L	L				B	3						3L							
	D	2		E	1		O	1				O	1	R	1											
2L	D	2	E	1	V	4	I	1	R	1	N	1	C	3	E	1	S	1	2W		2L					
	D	2			G	2		G	2	O	1	2L							2W							
	2W				N	1	3L		F	4		3L							2W							
P	3	I	1	L	1	I	1	S	1		T	1	U	1	T	1	O	1	R	1	I	1	A	1	L	1

1988

- Connect 4 solved, improved in 2015

1992

Gradient of NN w.r.t weights

- TD-Gammon

$$w_{t+1} - w_t = \alpha(Y_{t+1} - Y_t) \sum_{k=1}^t \lambda^{t-k} \nabla_w Y_k$$

Change weight from last turn

How important are older positions?

Change in board evaluation

1992



1992

64
6

TD Gammon

$$L = \alpha ||V_t(s_{t+1}) - V_t(s_t)||$$

We want predictions about present to match future.

TD Gammon

$$L = \alpha ||V_t(s_{t+1}) - V_t(s_t)||$$

We want predictions about present to match future.

$$L = \alpha ||z - V_t(s_t)||$$

Z is the true outcome.

TD Gammon

$$L = \alpha ||V_t(s_{t+1}) - V_t(s_t)||$$

Leads to parameter updates: δ_t is difference between two state predictions.

$$\theta_{t+1} = \theta_t + \alpha \delta_t \frac{\partial V_t}{\partial \theta_t}$$

But. Need to assign credit.

TD Gammon

$$\theta_{t+1} = \theta_t + \alpha \delta_d \frac{\partial V_t}{\partial \theta_t}$$

But. Need to assign credit.
Store ALL gradients. Weight:

$$\theta_{t+1} = \theta_t + \alpha \delta_t \sum_{k=1}^t \lambda^{t-k} \frac{\partial V_k}{\partial \theta_t}$$

1996-7

- Deep Blue loses to, then beats Kasparov



2007

- Draughts solved, search space: 5×10^{20}

RL in Games: State of the Art

Program	Level of Play	RL Program to Achieve Level
Checkers	Perfect	<i>Chinook</i>
Chess	International Master	<i>KnightCap / Meep</i>
Othello	Superhuman	<i>Logistello</i>
Backgammon	Superhuman	<i>TD-Gammon</i>
Scrabble	Superhuman	<i>Maven</i>
Go	Grandmaster	<i>MoGo¹, Crazy Stone², Zen³</i>
Poker ⁴	Superhuman	<i>SmooCT</i>

¹9 × 9²9 × 9 and 19 × 19³19 × 19⁴Heads-up Limit Texas Hold'em

2015

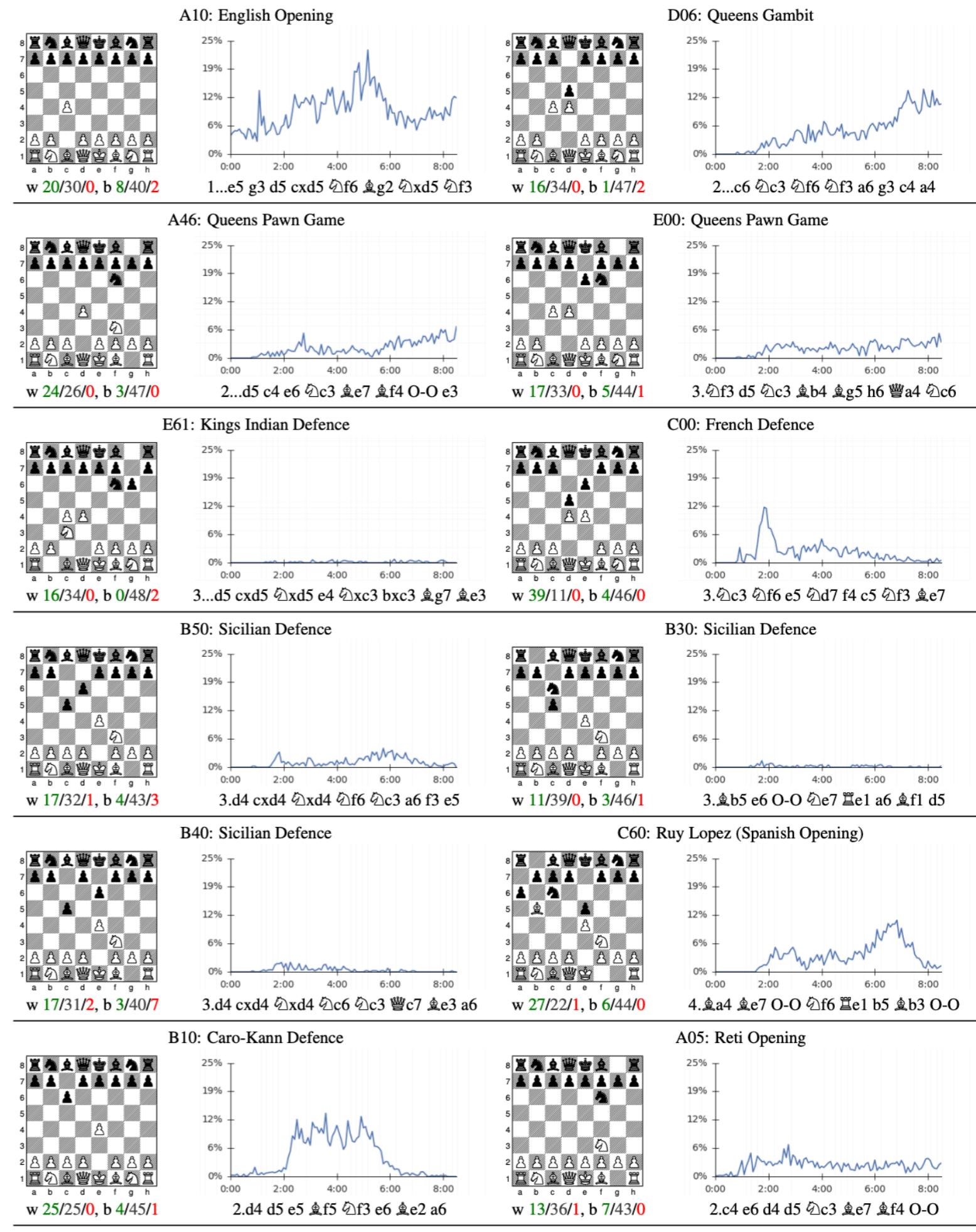
- AlphaGo

2017

- AlphaGo Zero beats AlphaGo 100-0

2017

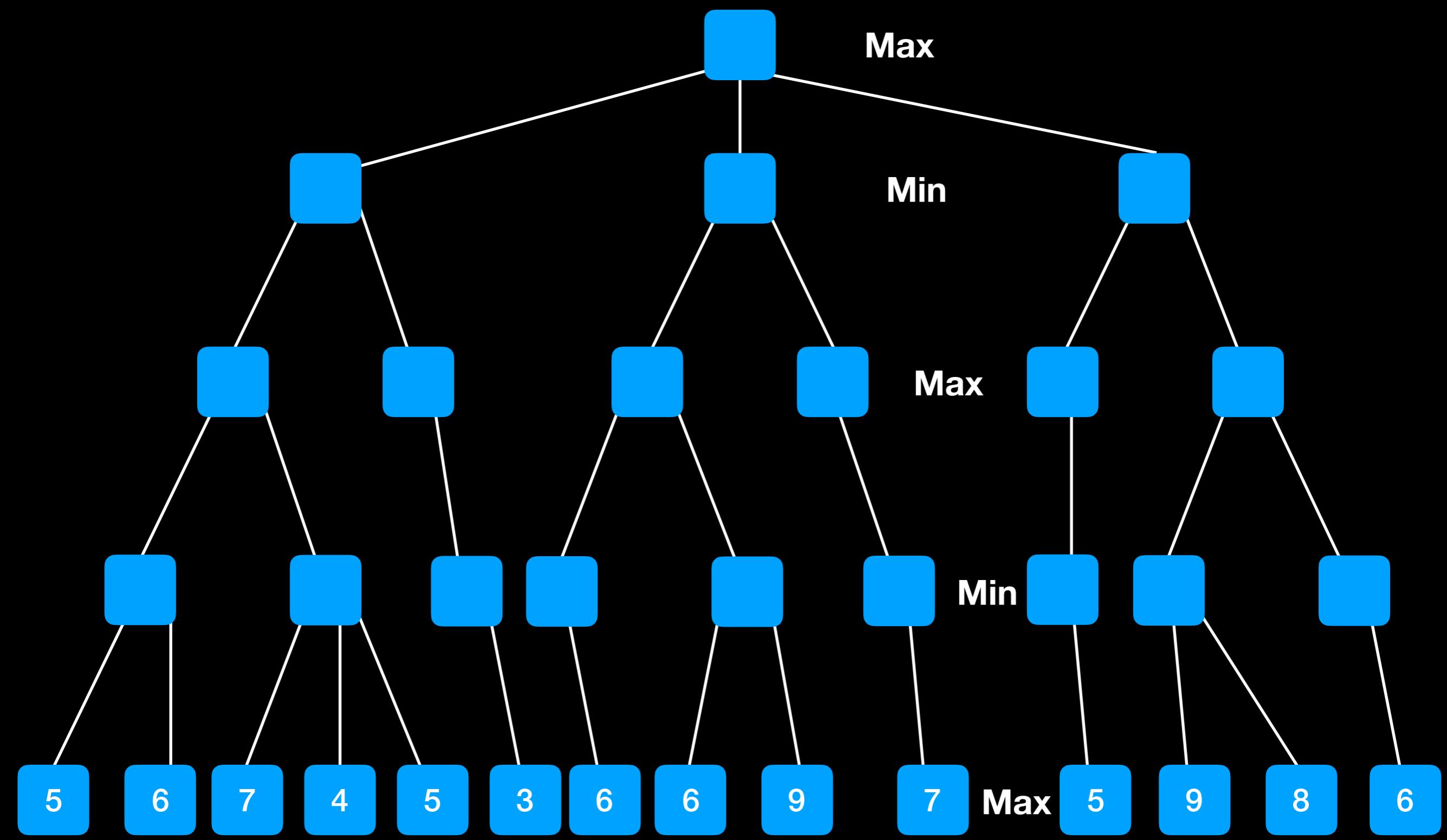
- AlphaZero beats a version of Stockfish 28-0-78

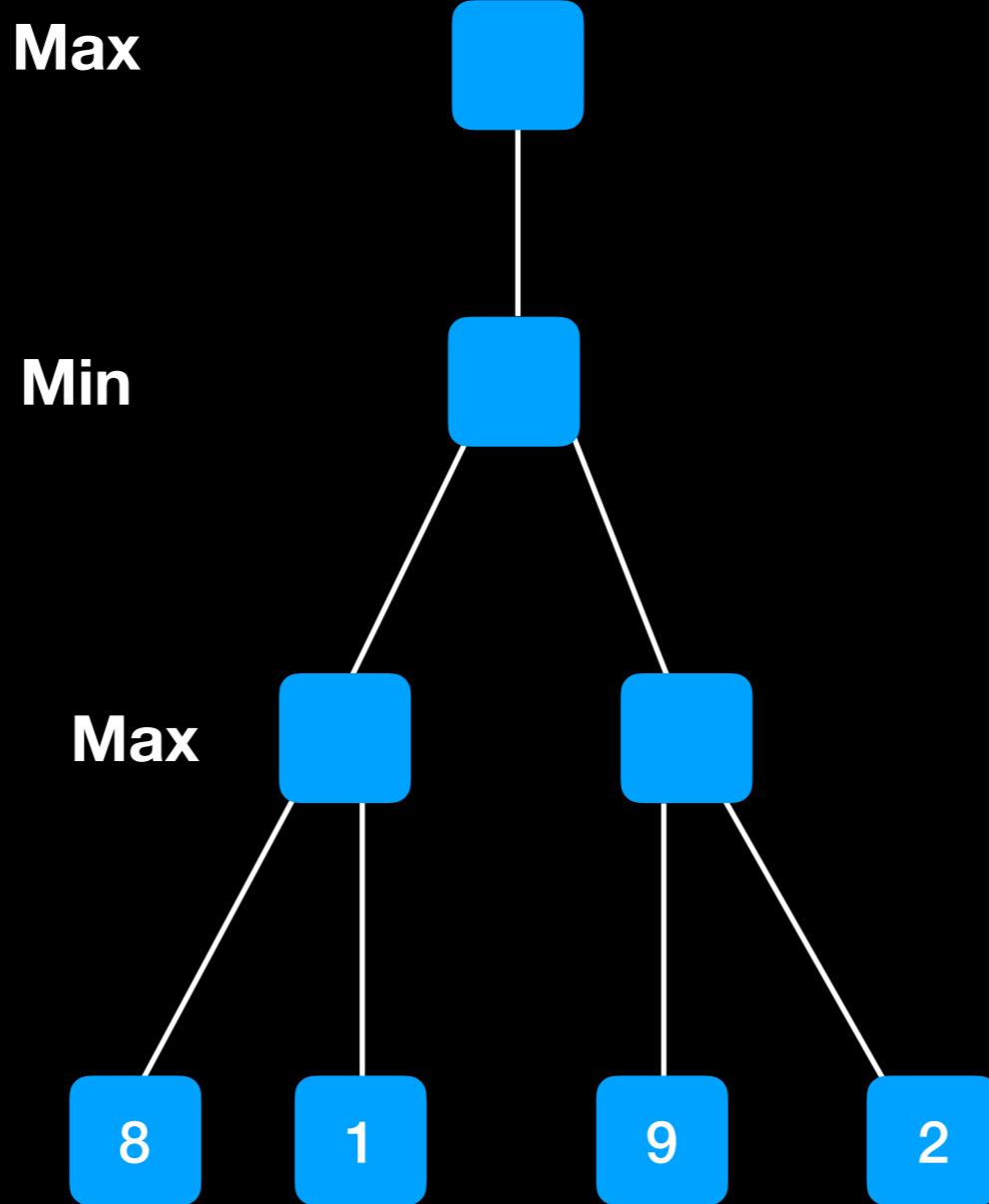


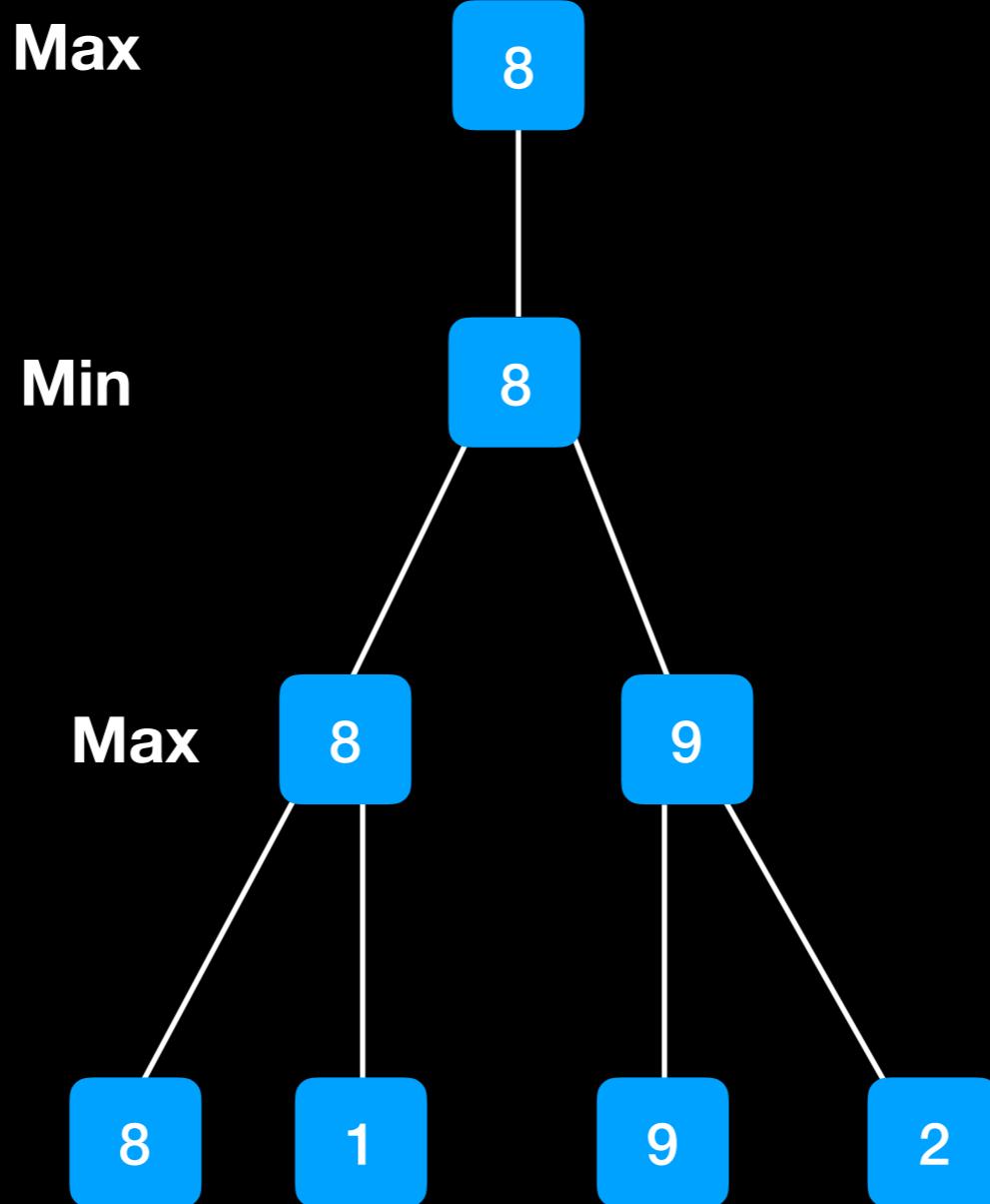
2018

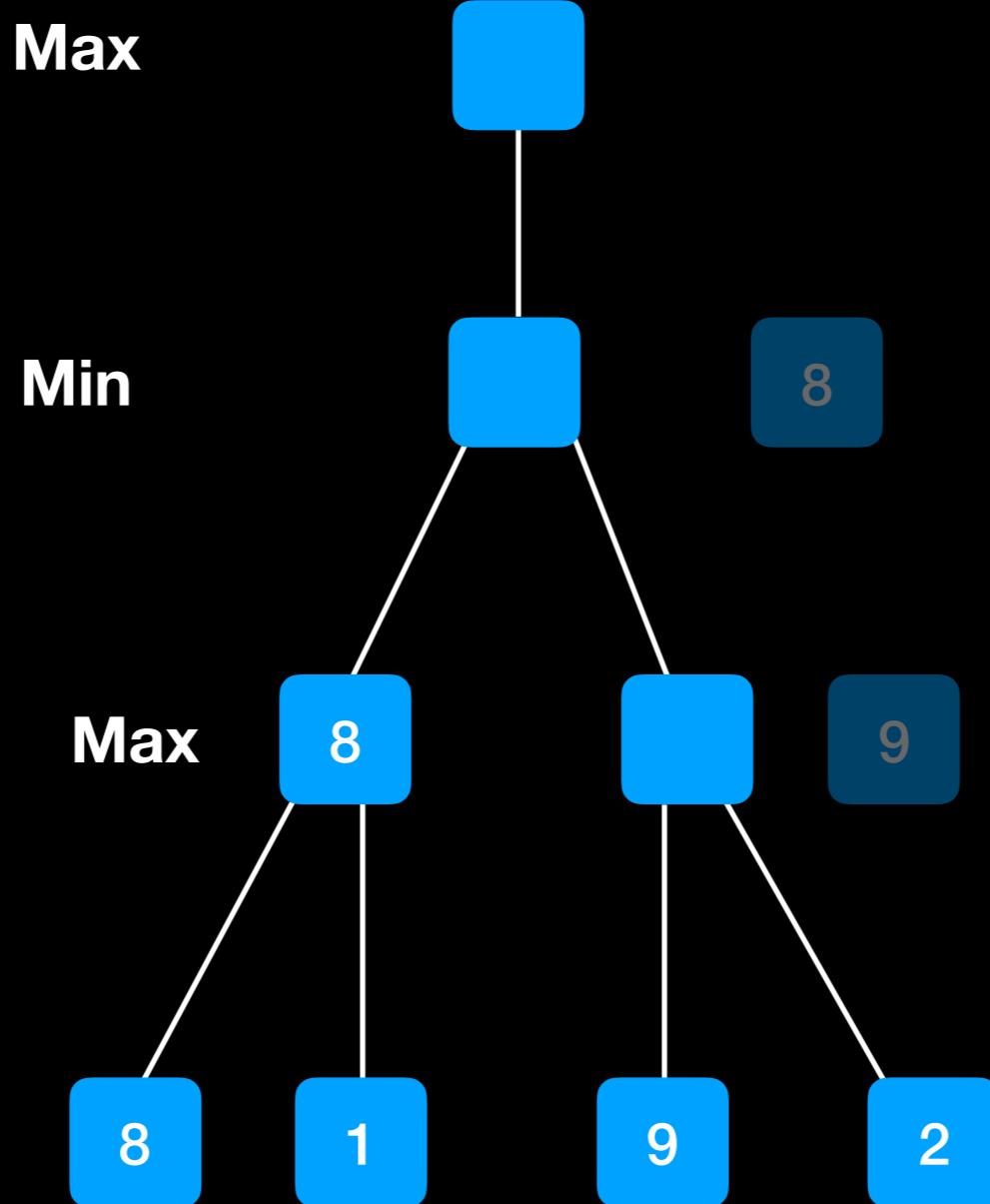


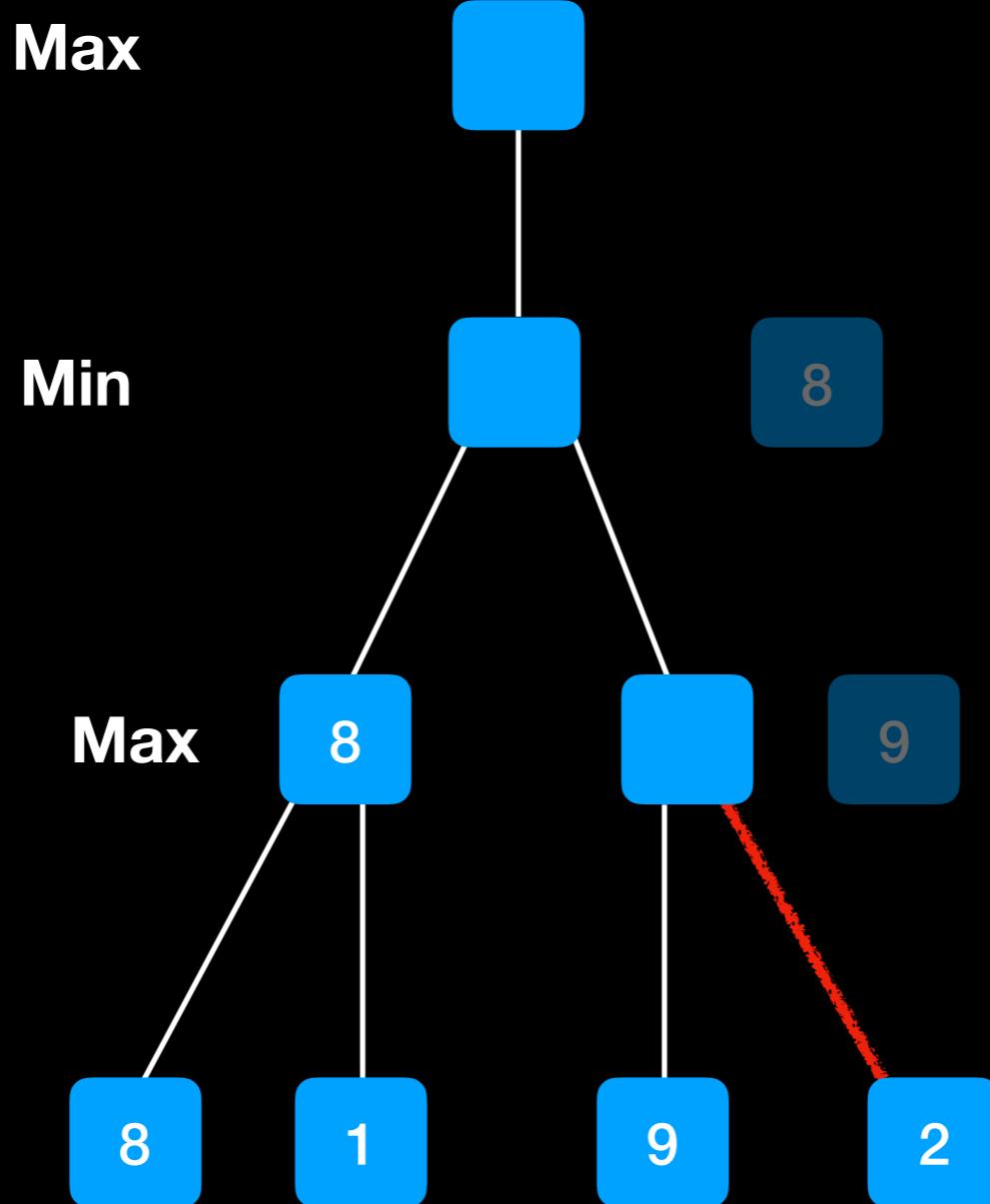




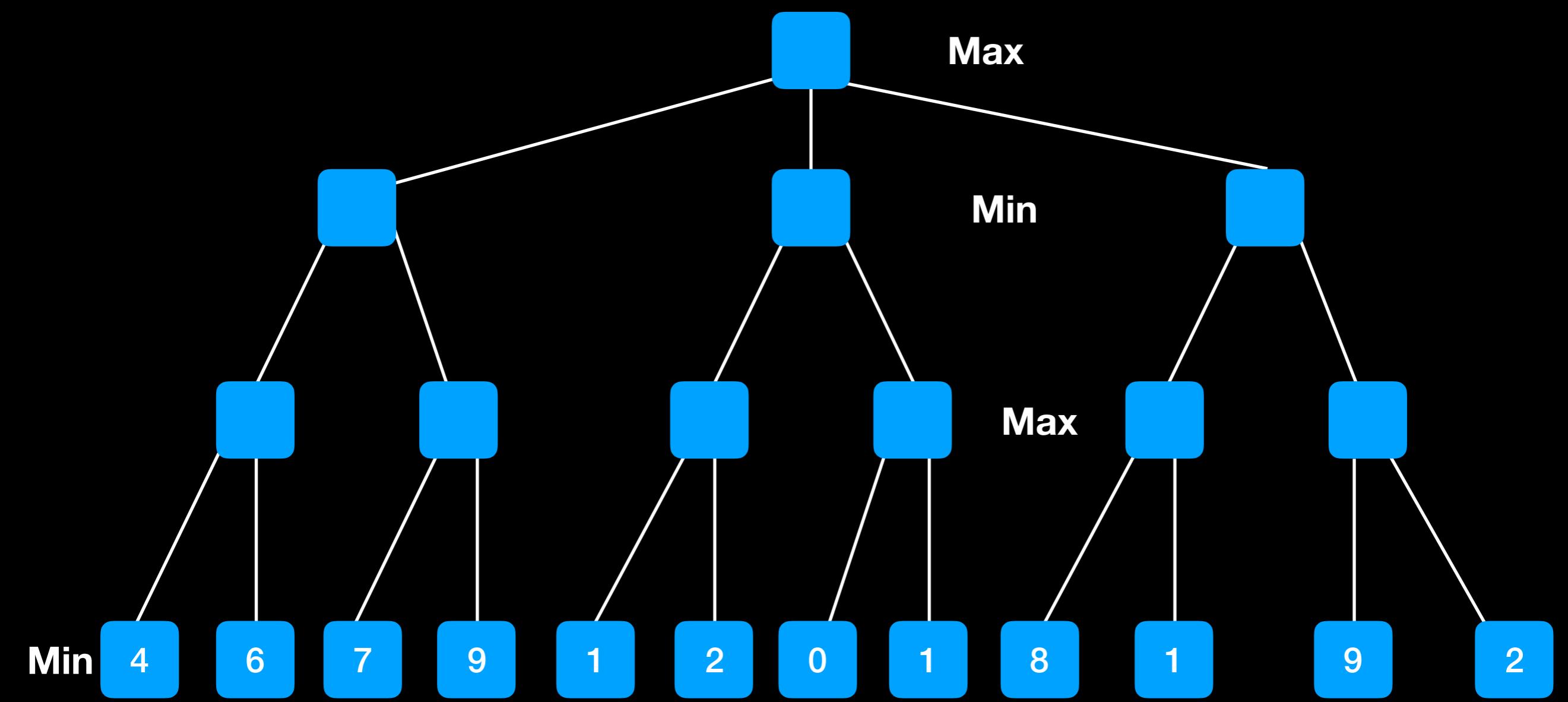


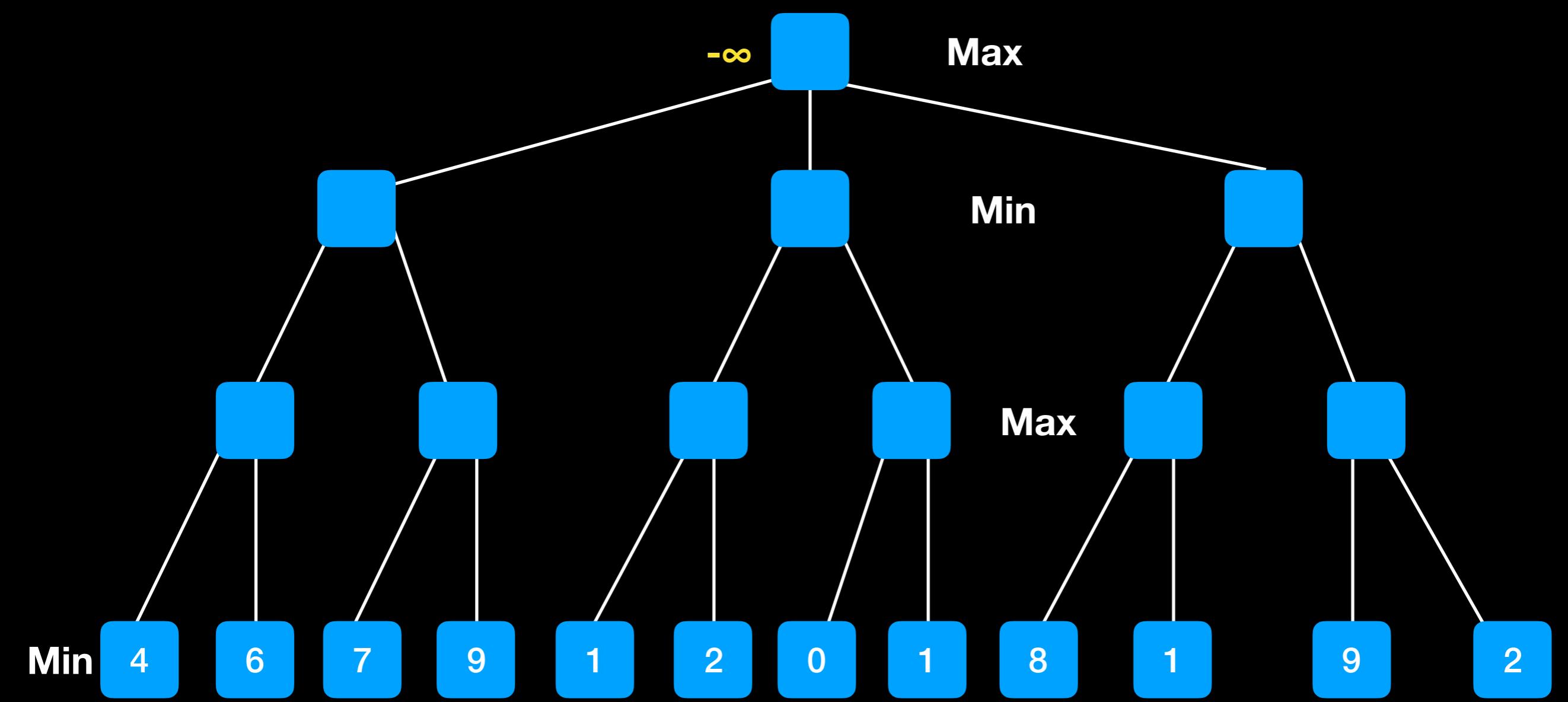


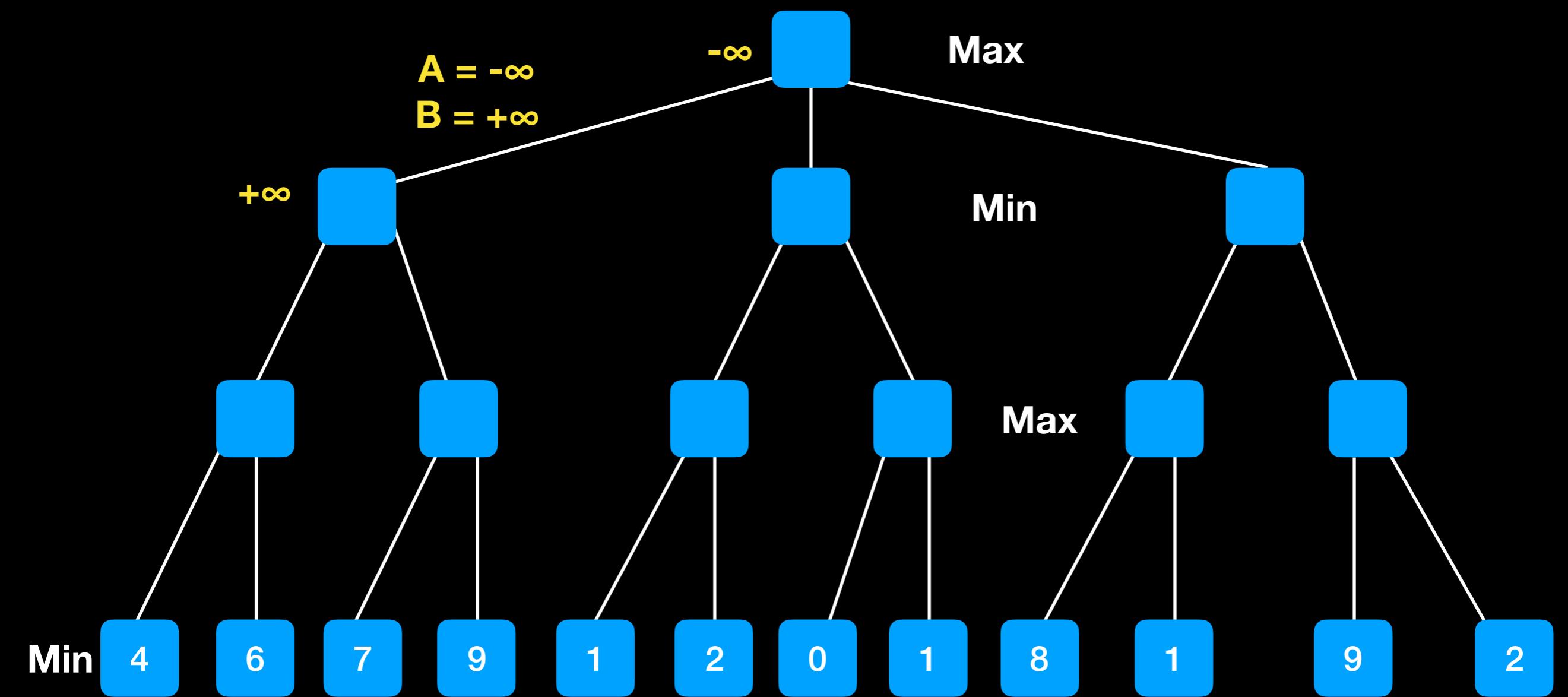


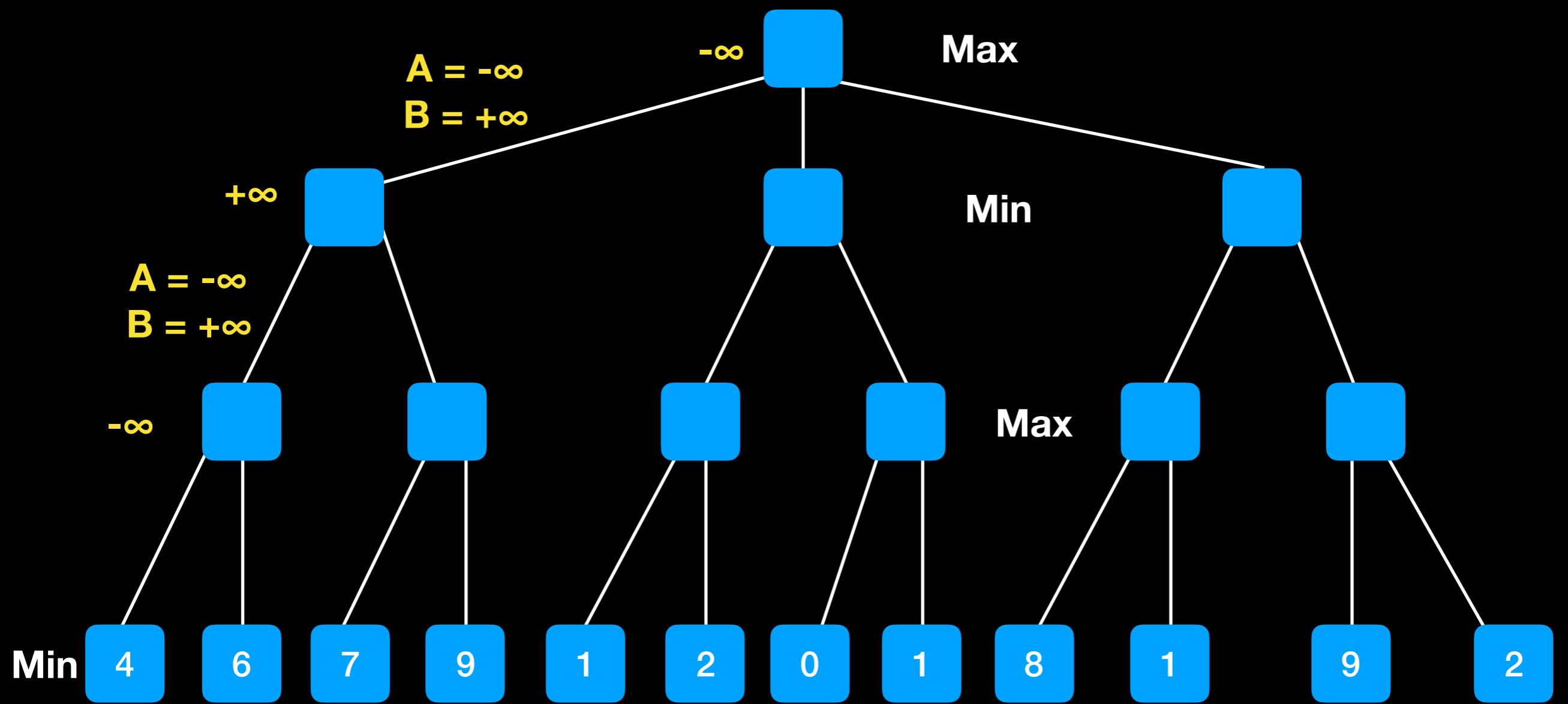


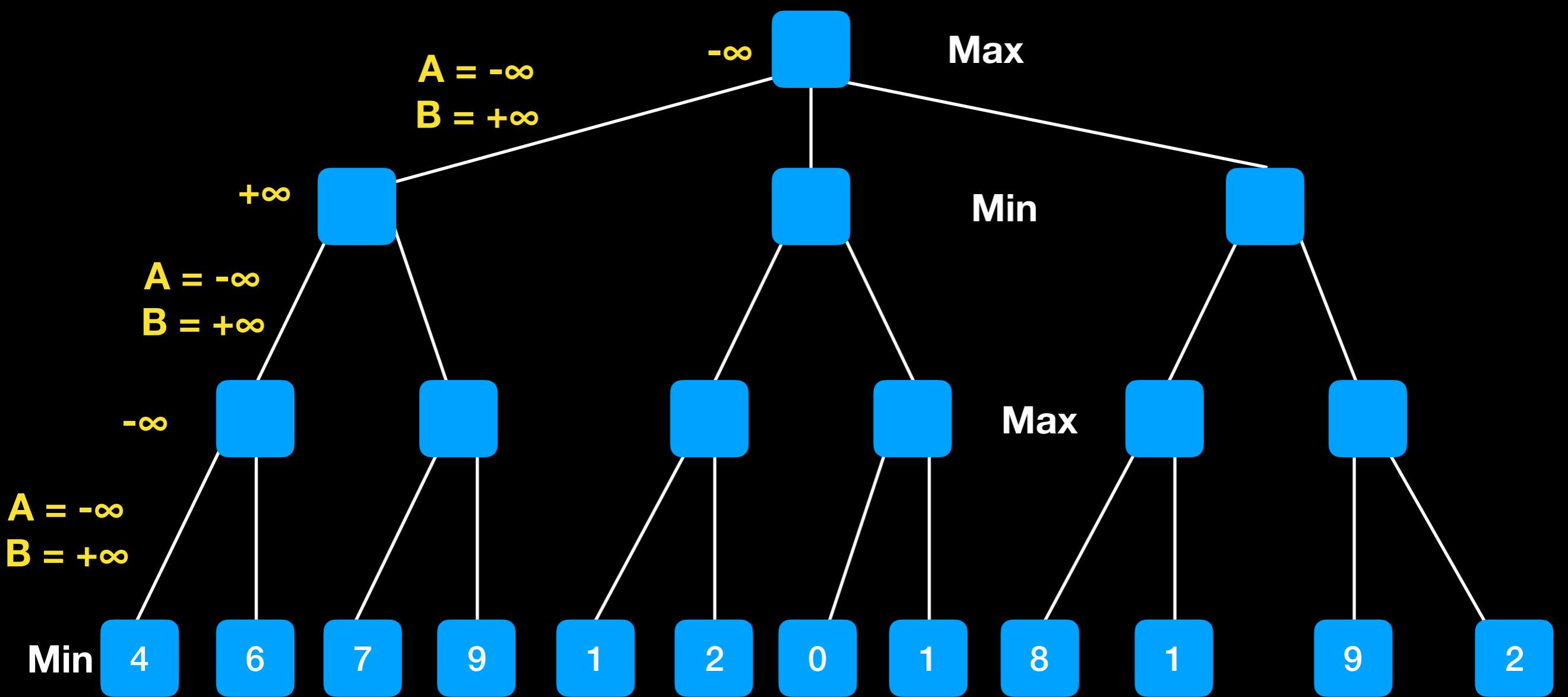
```
function alphabeta(node, depth, α, β, maximizingPlayer) is
    if depth = 0 or node is a terminal node then
        return the heuristic value of node
    if maximizingPlayer then
        value := -∞
        for each child of node do
            value := max(value, alphabeta(child, depth - 1, α, β, FALSE))
            α := max(α, value)
            if α ≥ β then
                break (* β cut-off *)
        return value
    else
        value := +∞
        for each child of node do
            value := min(value, alphabeta(child, depth - 1, α, β, TRUE))
            β := min(β, value)
            if α ≥ β then
                break (* α cut-off *)
        return value
```

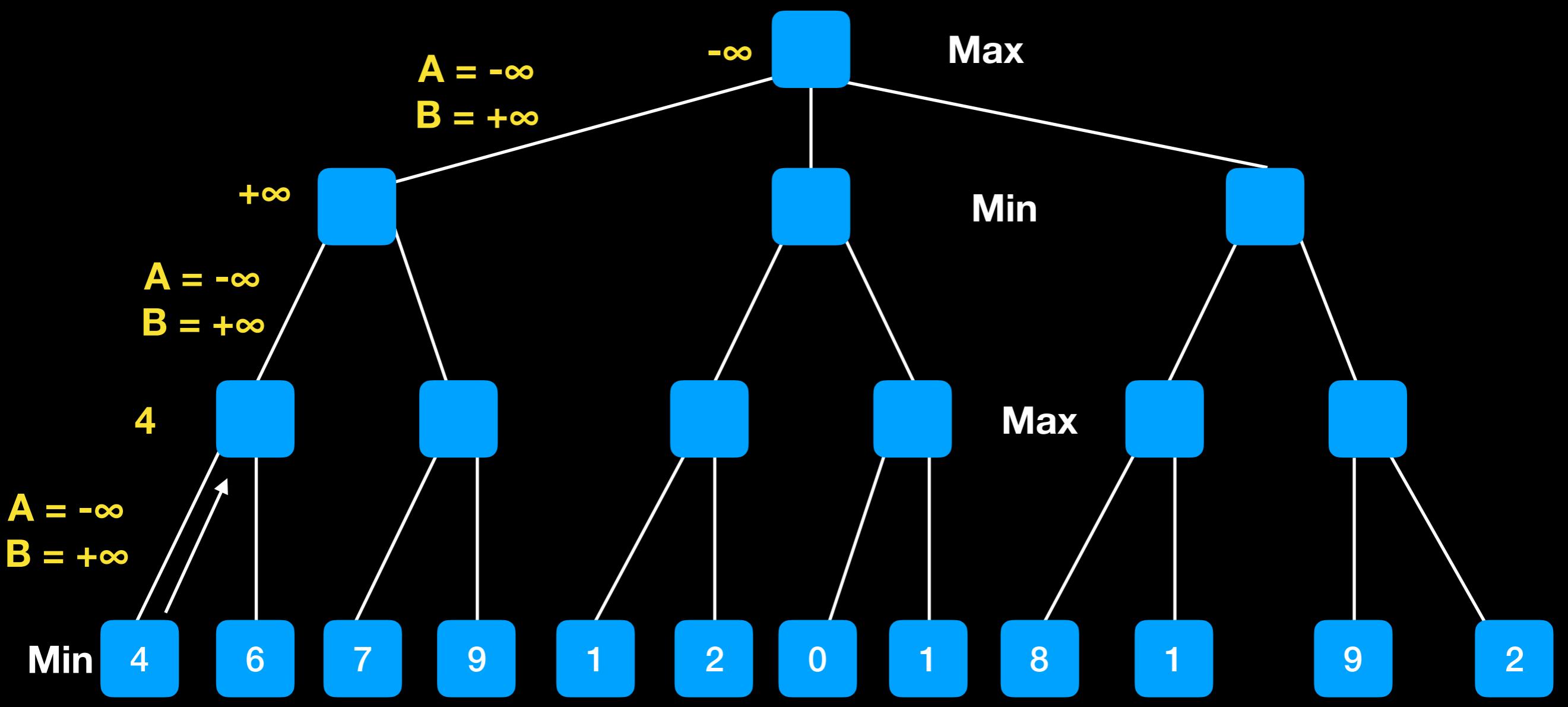




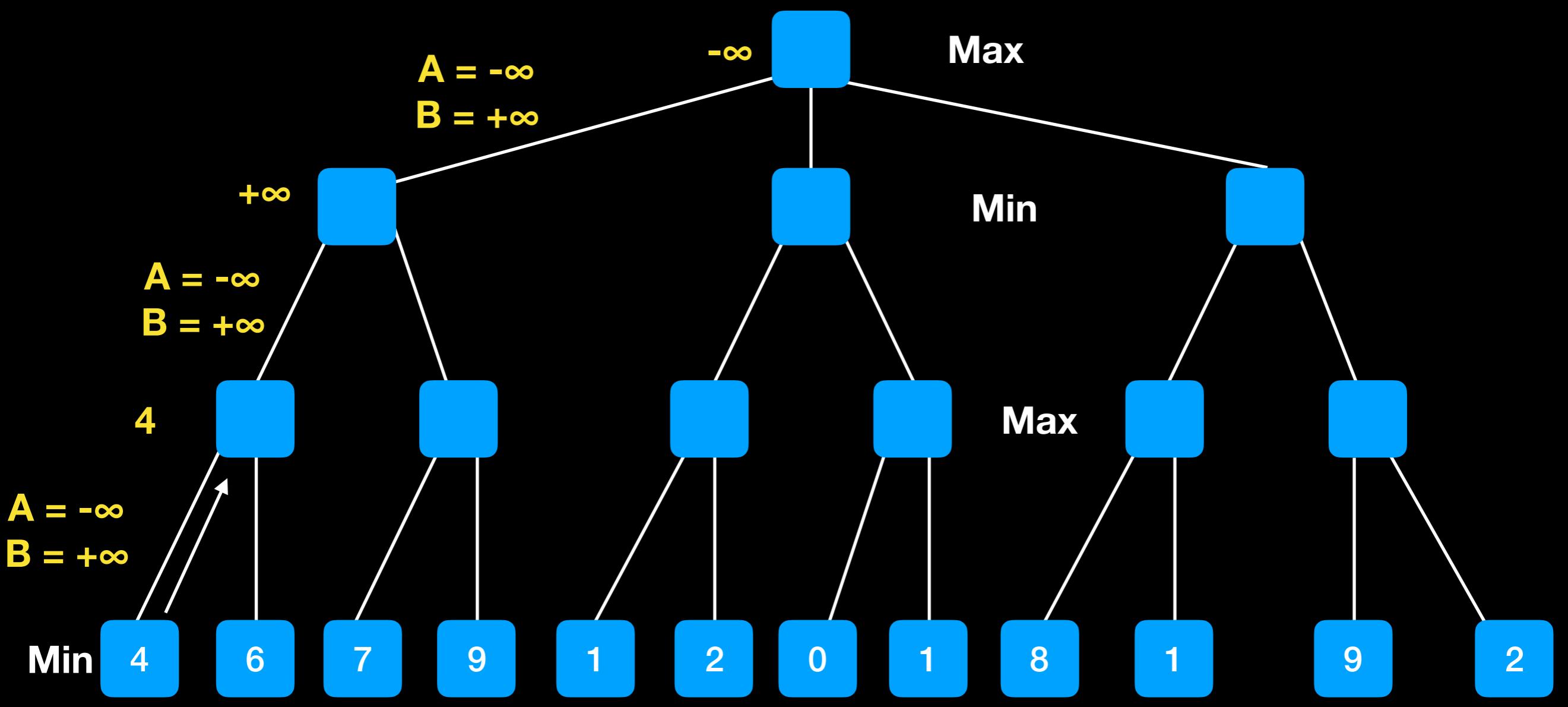




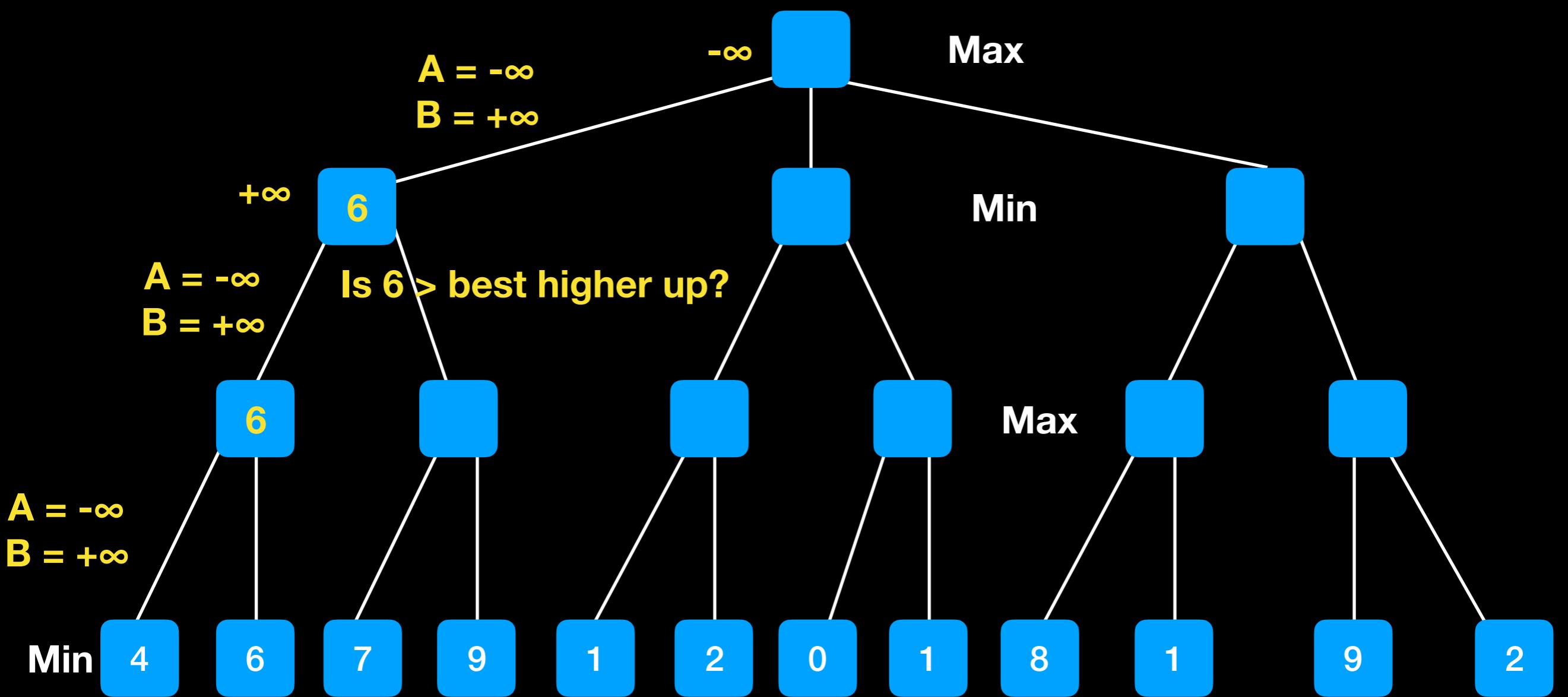


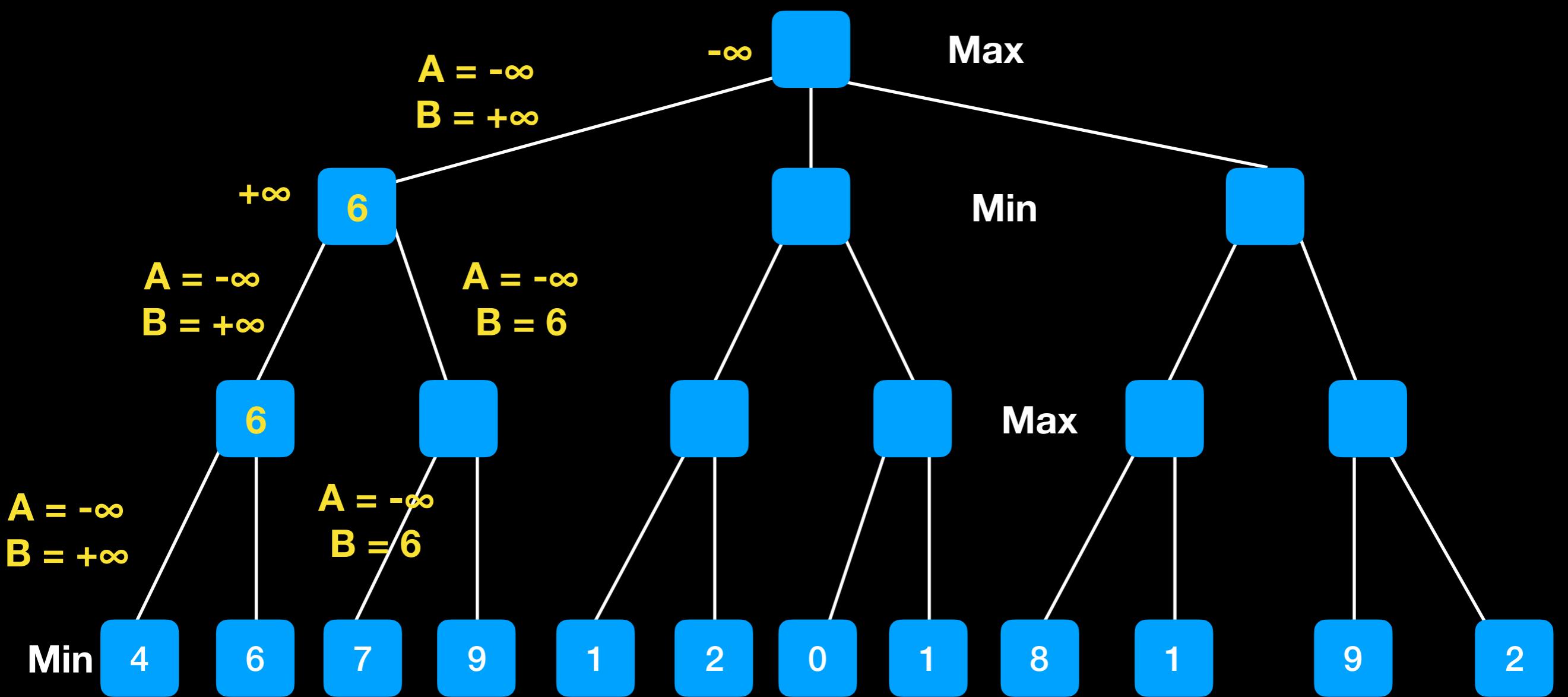


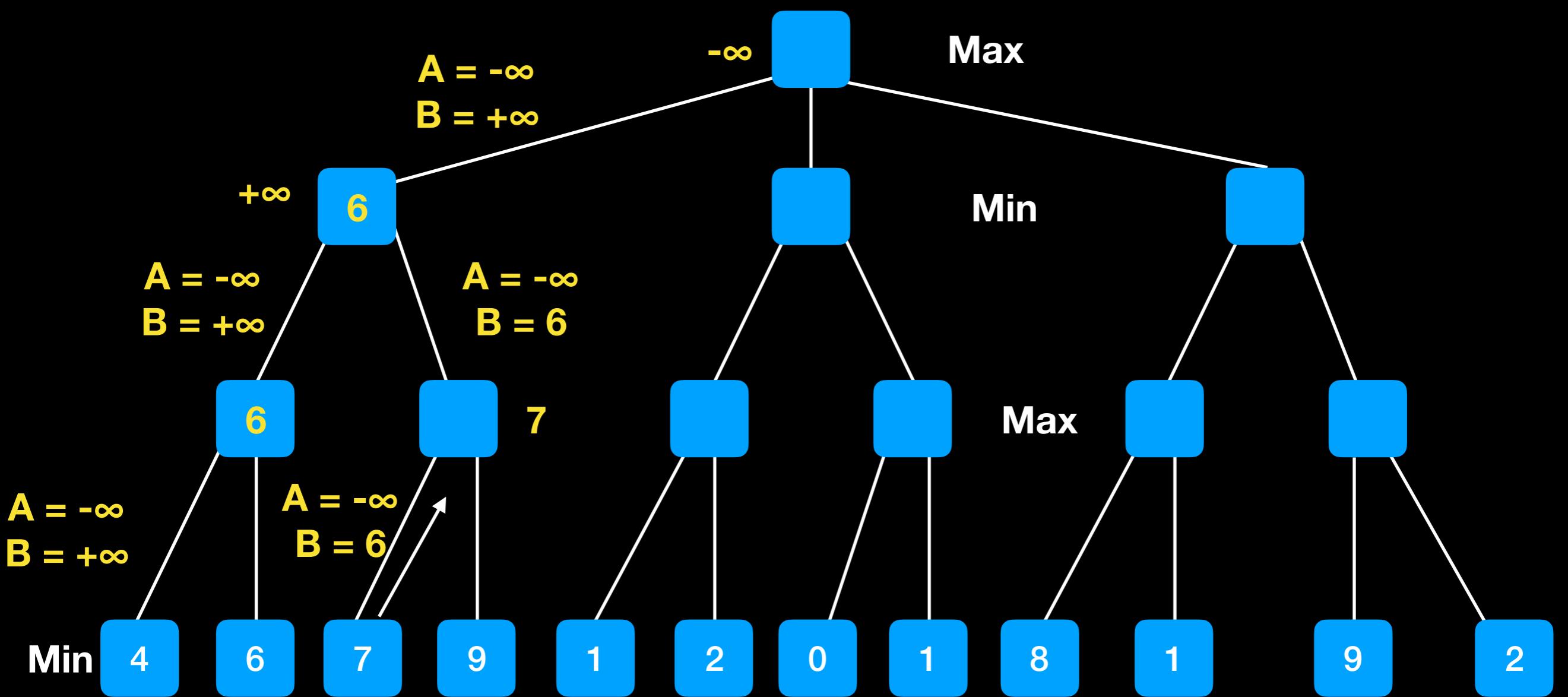
Is 4 > B?

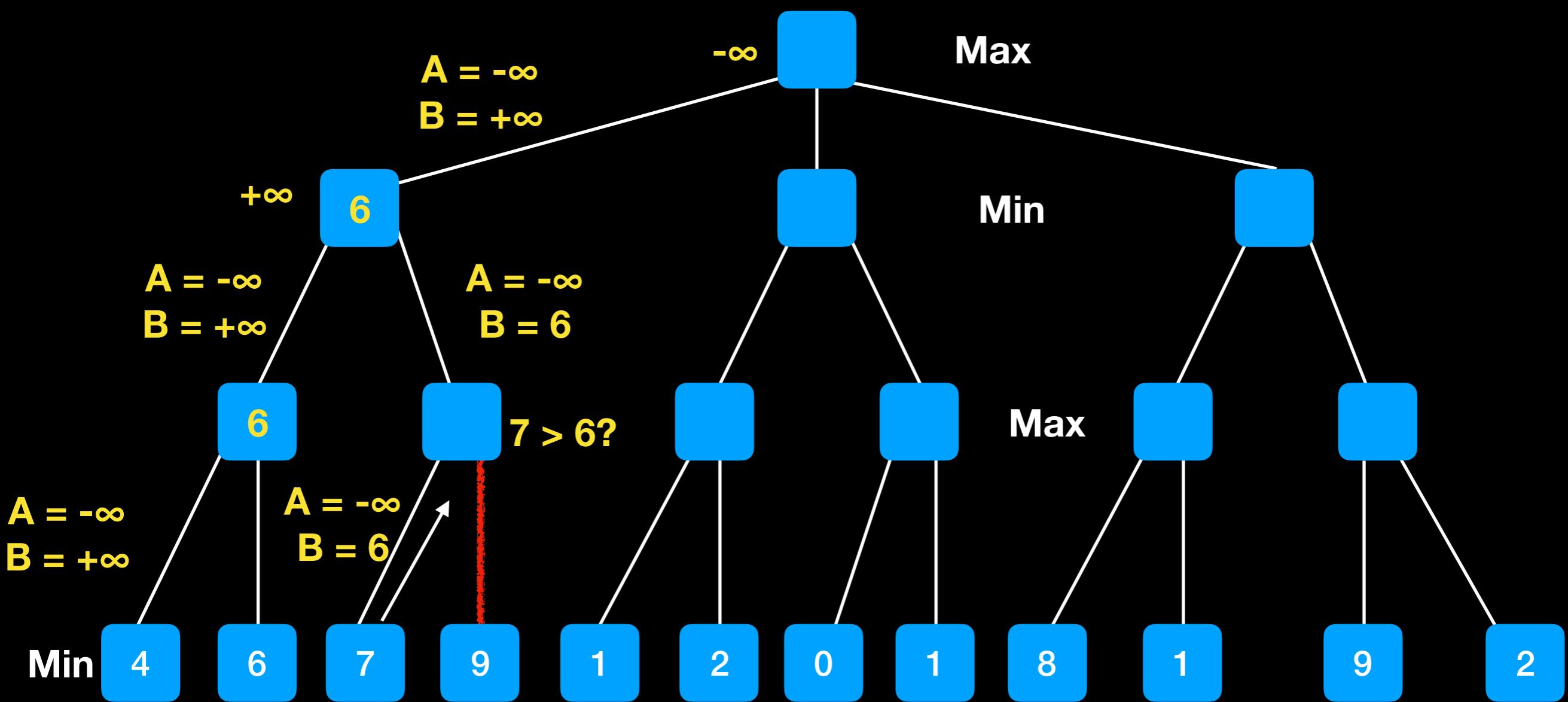


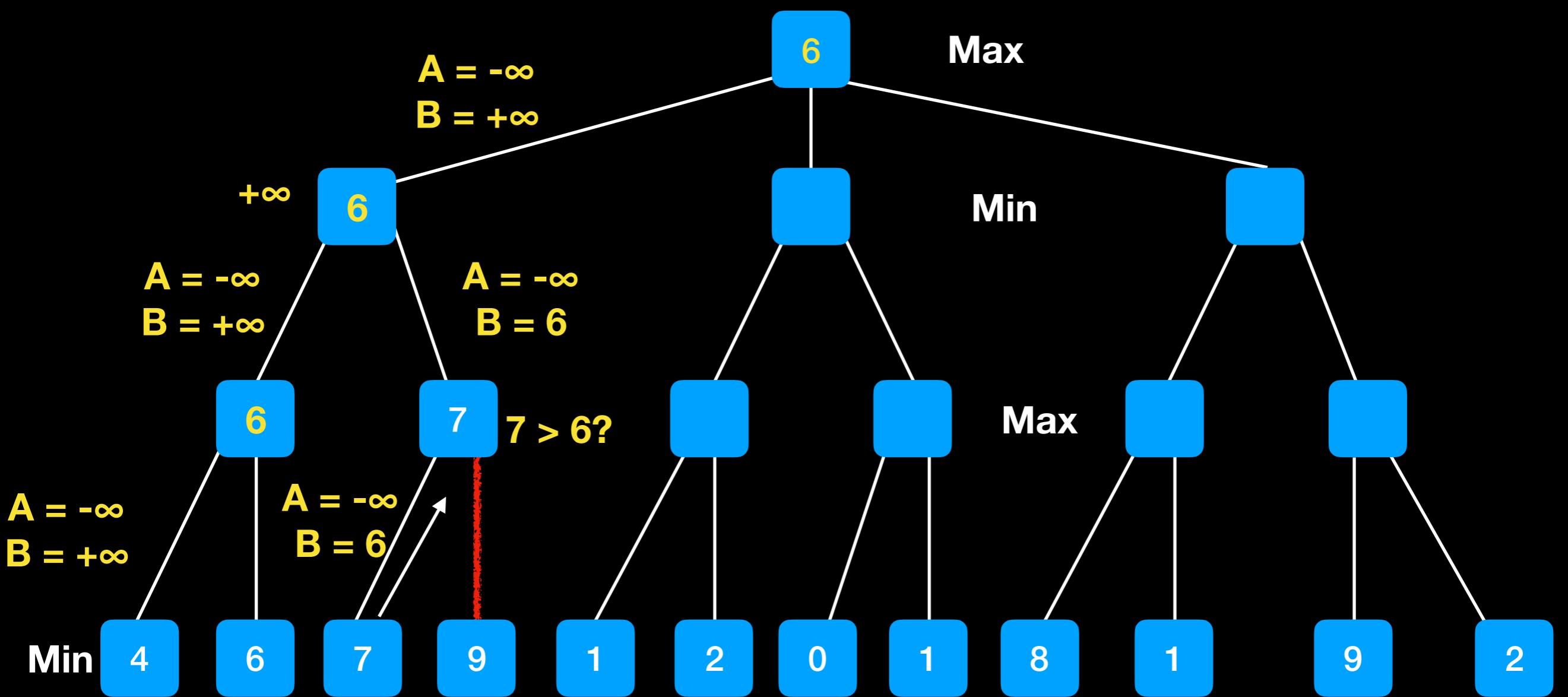
Is 4 > B?

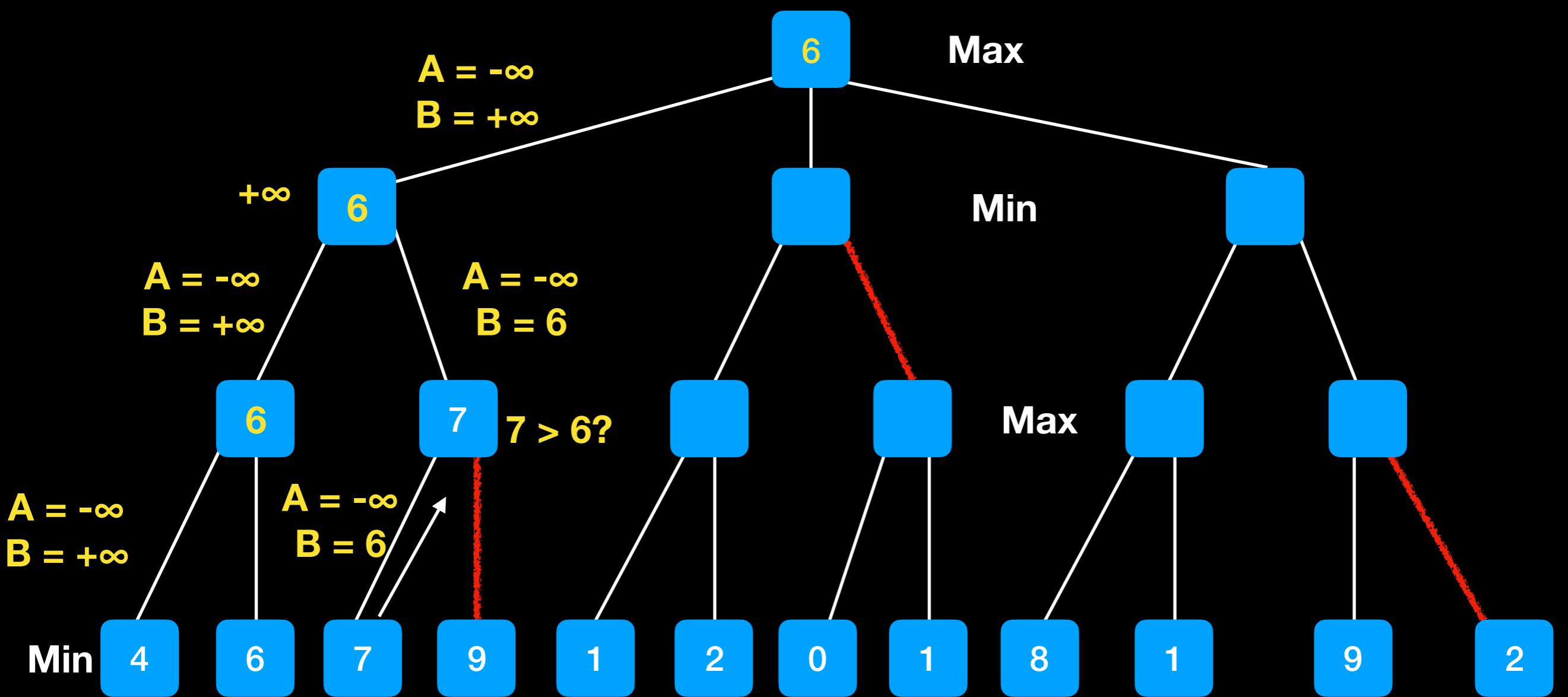


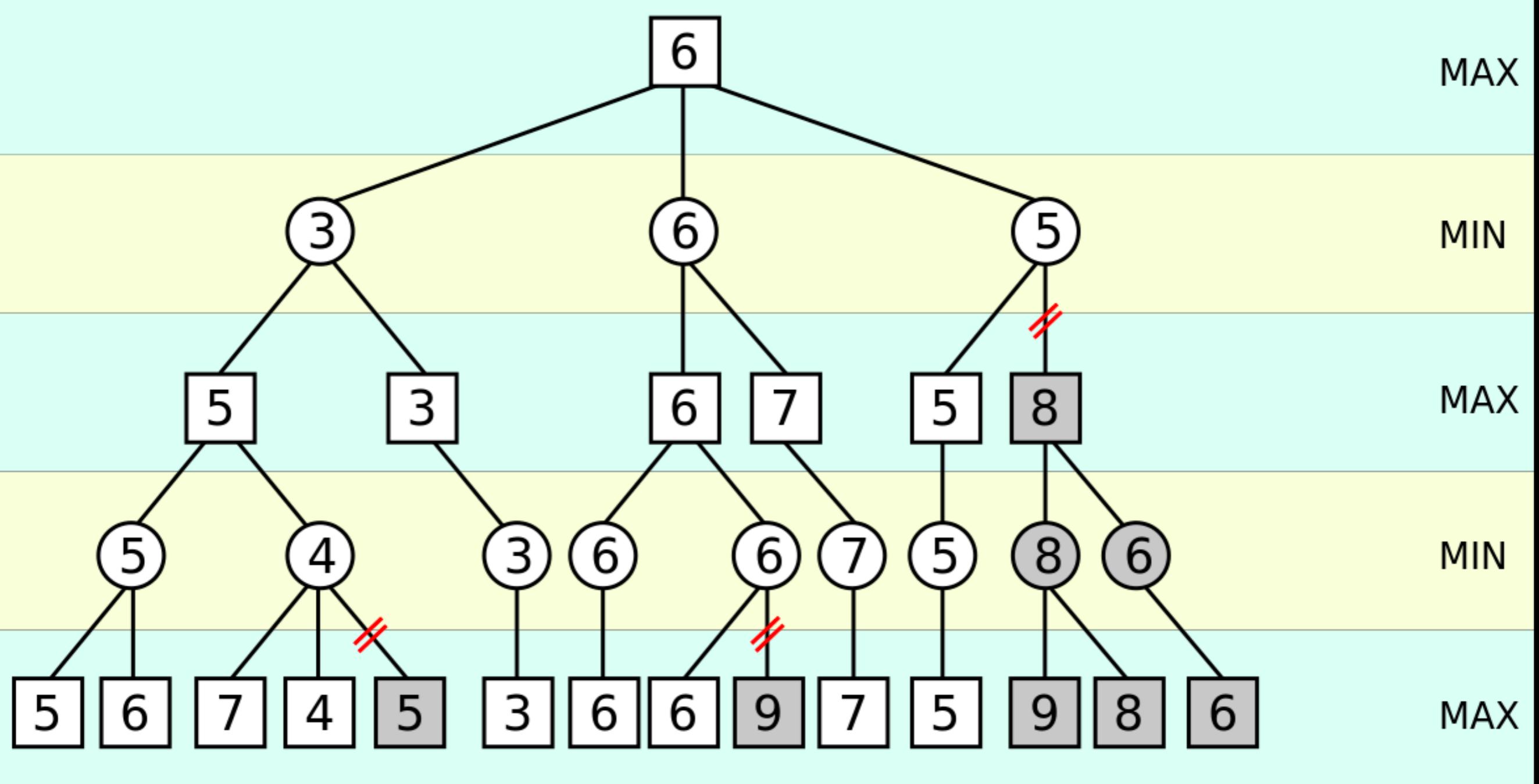












No. of moves = b^d

No. of moves = b^d

chess: $b \sim 35, d \sim 80$

No. of moves = b^d

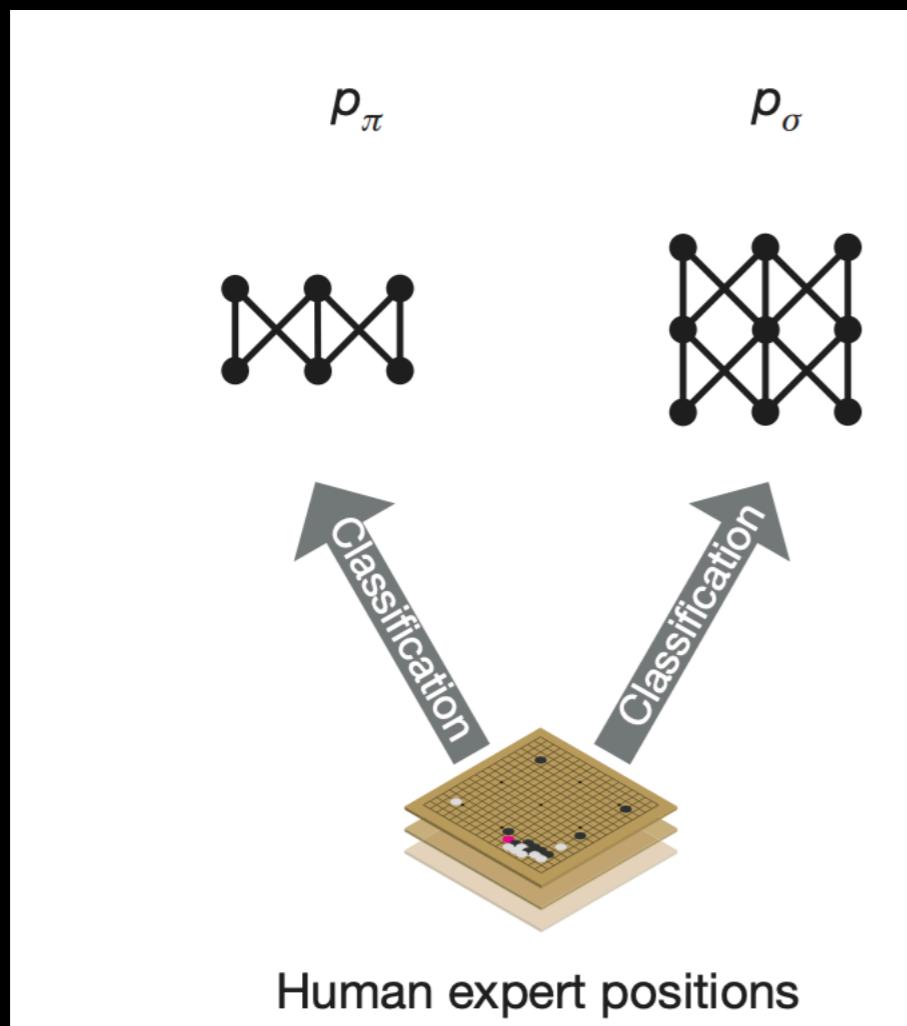
chess: $b \sim 35, d \sim 80$

Go: $b \sim 250, d \sim 150$



Rollout Policy *Fast*

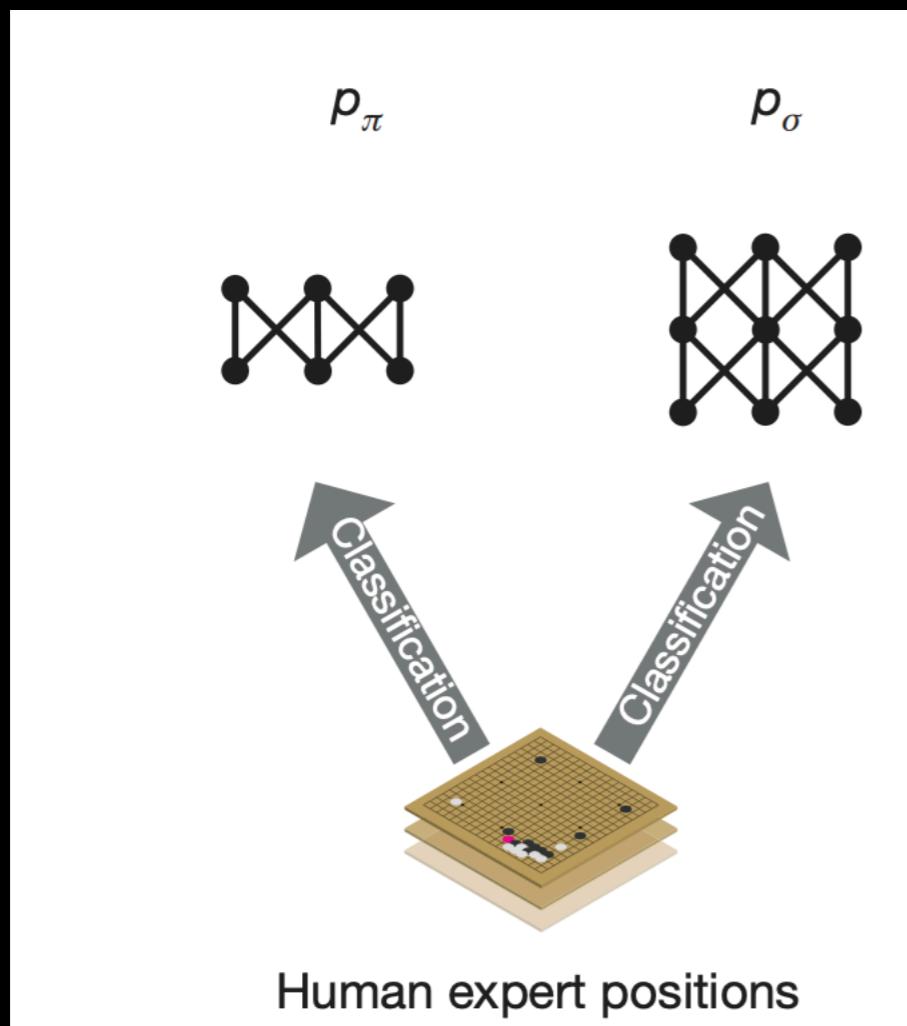
SL Policy



Used online game database of top players
30 million positions.

Rollout Policy
Fast

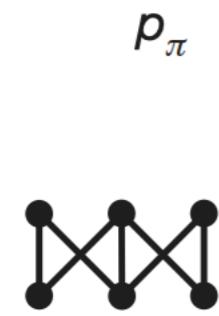
SL Policy



Used online game database of top players
30 million positions.

GOAL: Find human master move.

Rollout Policy
Fast



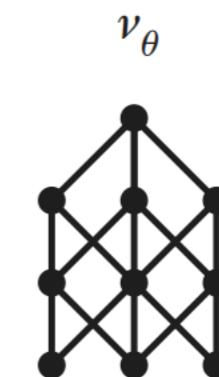
SL Policy



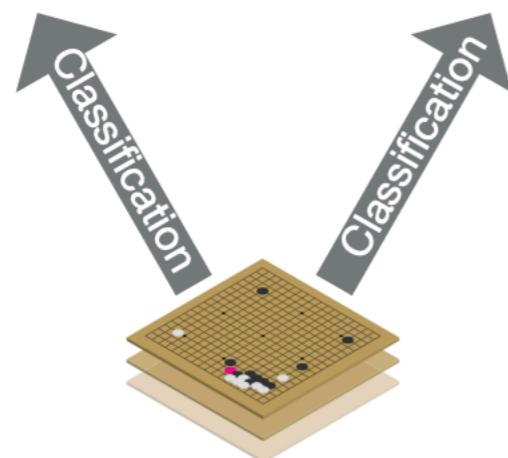
RL Policy Network



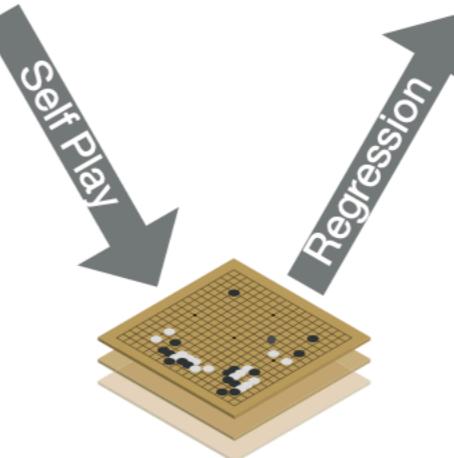
Value Network
Predict winner



Policy gradient →

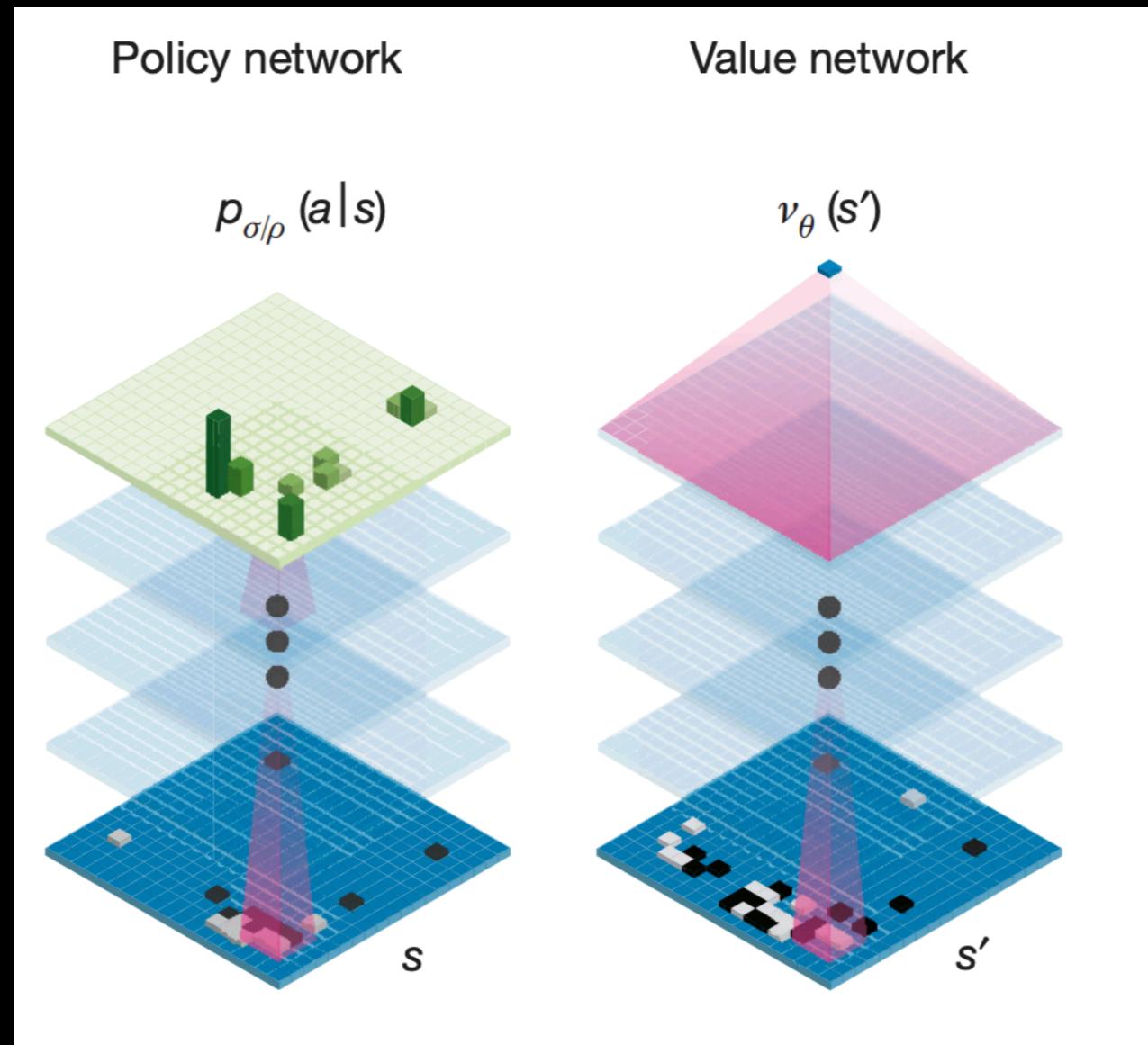


Human expert positions



Self-play positions

GOAL: Win



Monte Carlo Tree Search

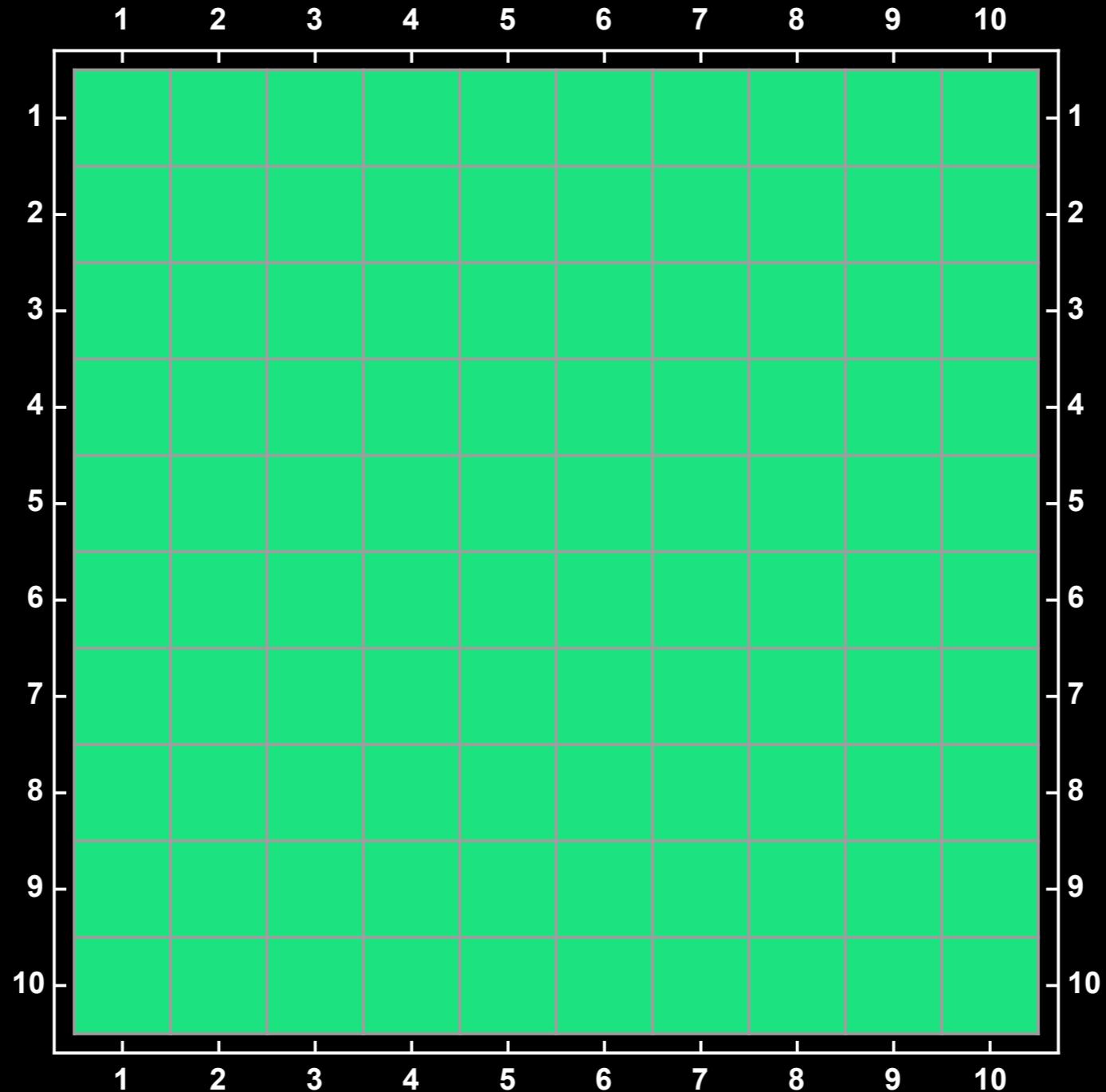
Monte Carlo Tree Search

Random Rollout

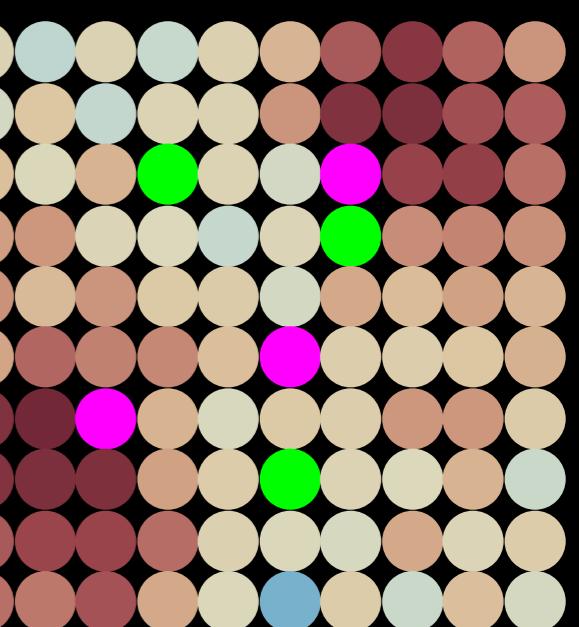
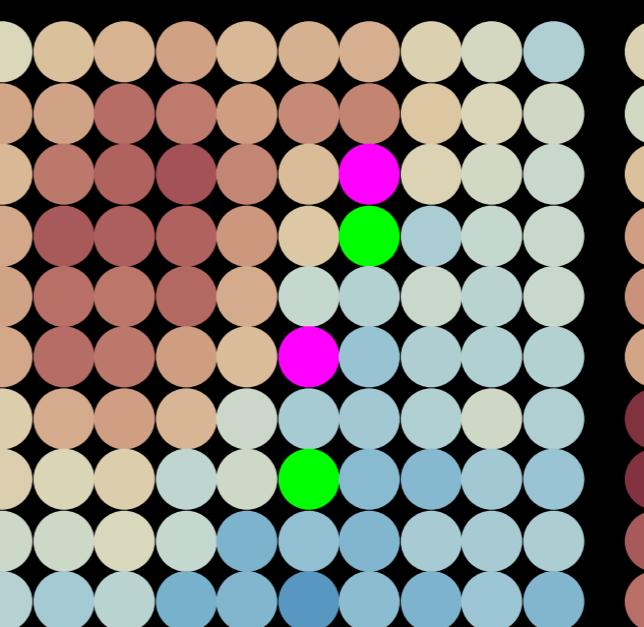
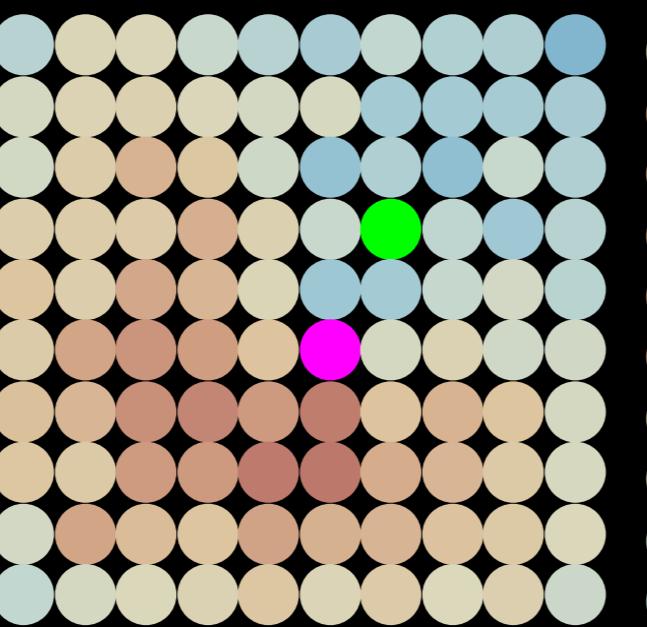
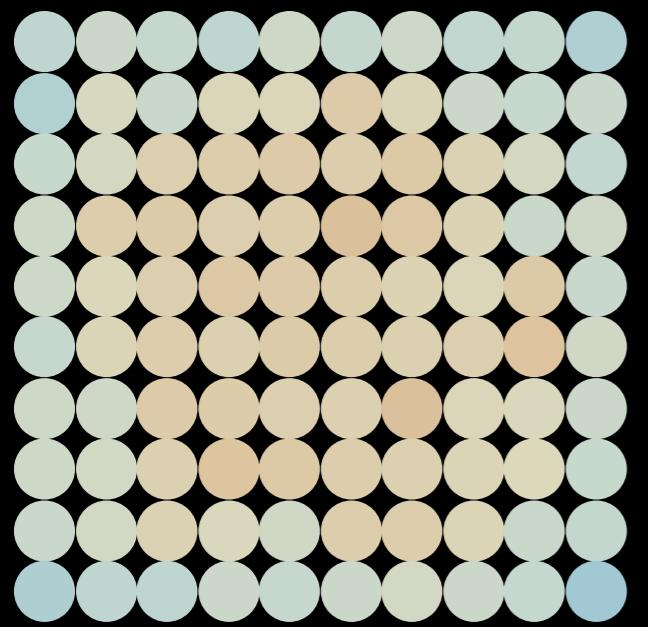
$$V(s_L) = (1 - \lambda)v_\theta(s_L) + \lambda z_L$$

Value Network

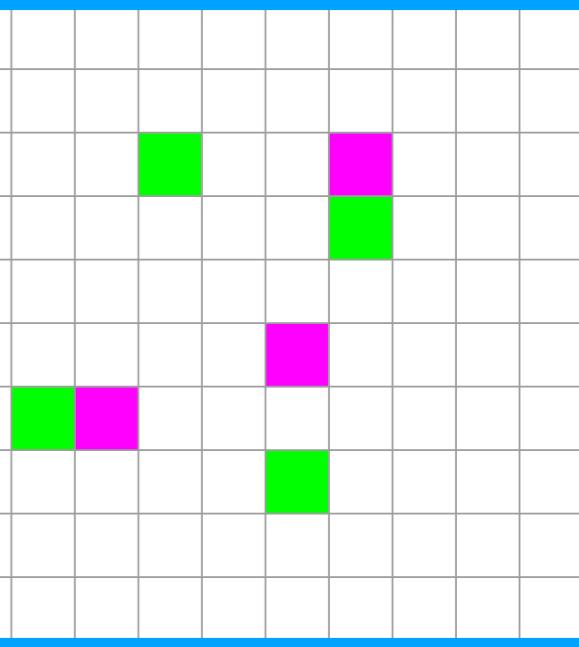
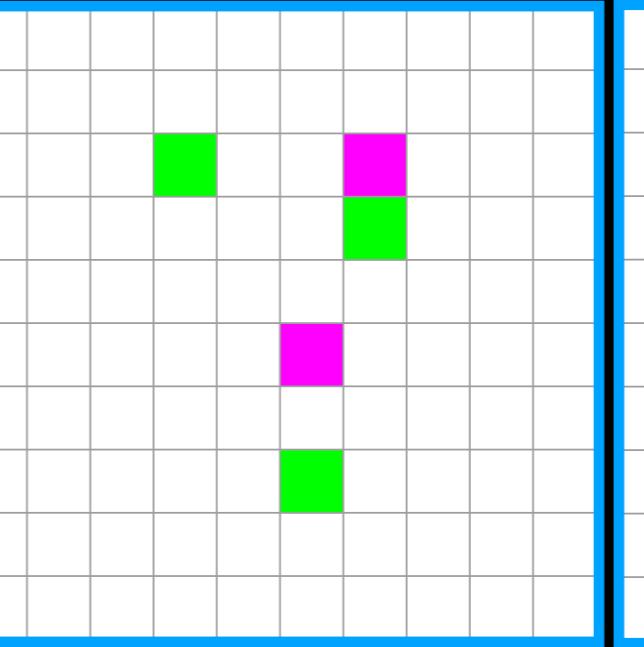
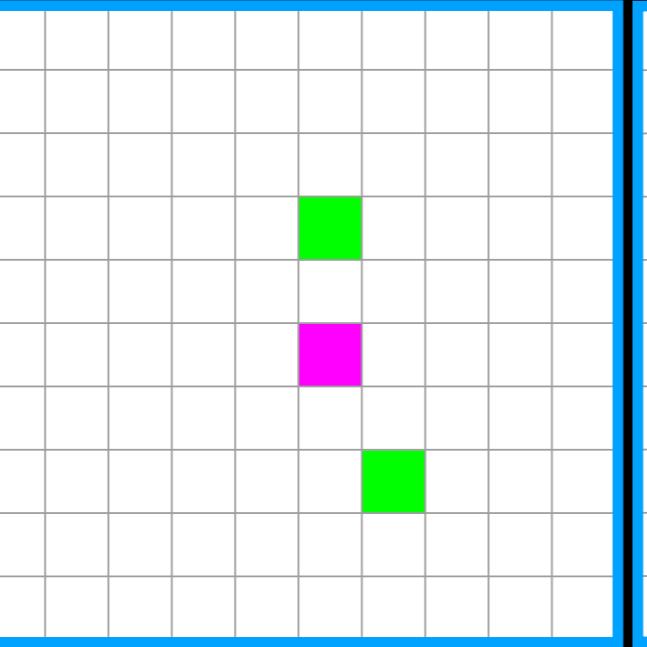
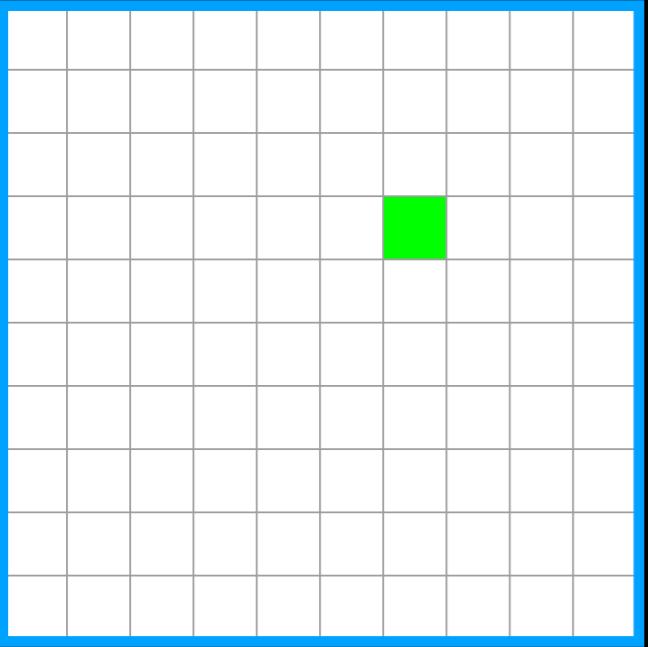
DEMO



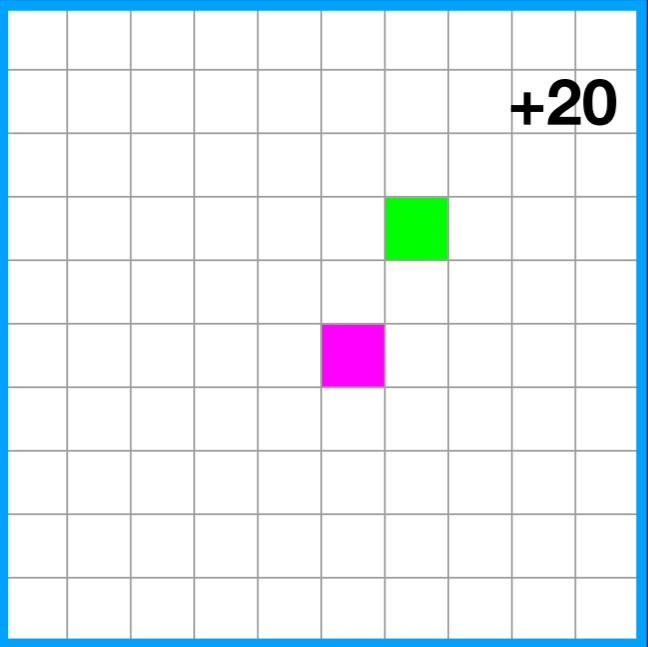
Probability



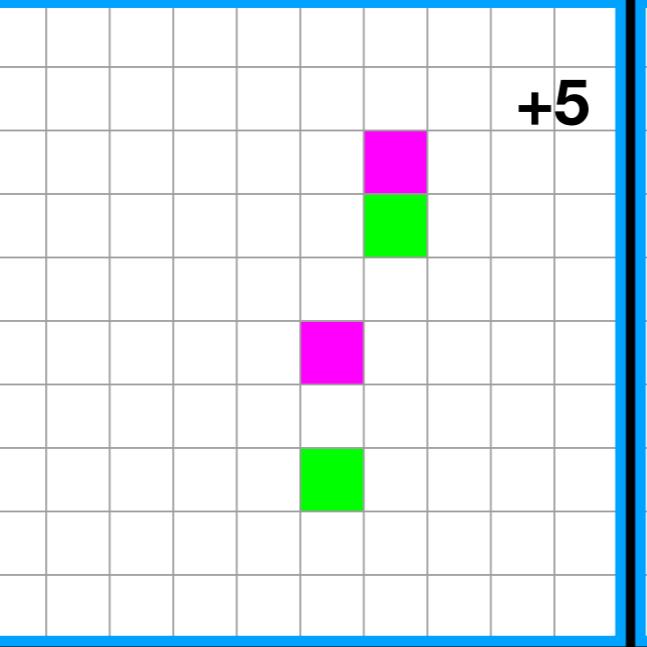
Computer Move



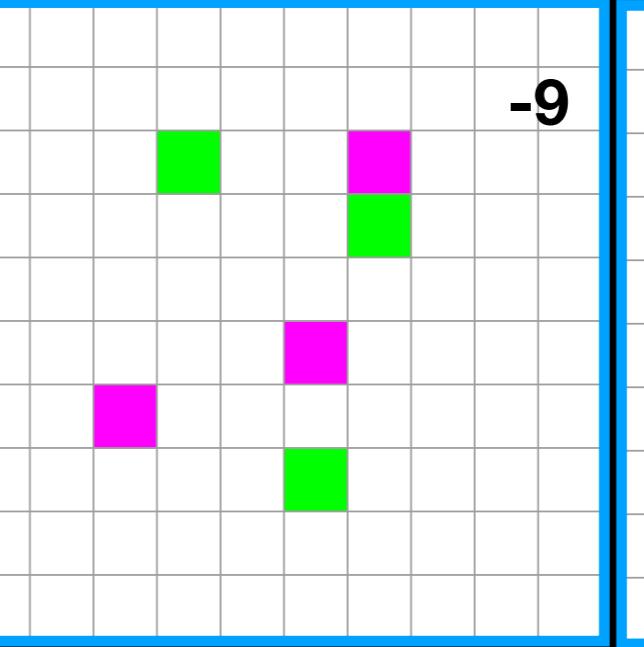
My Move



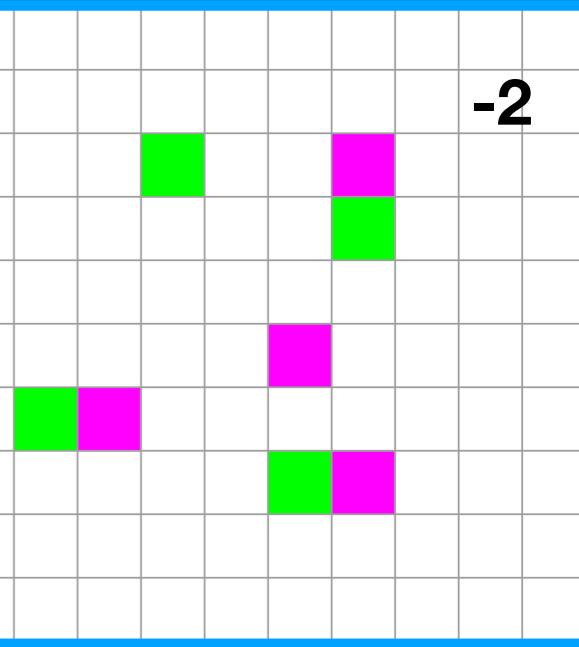
+20



+5

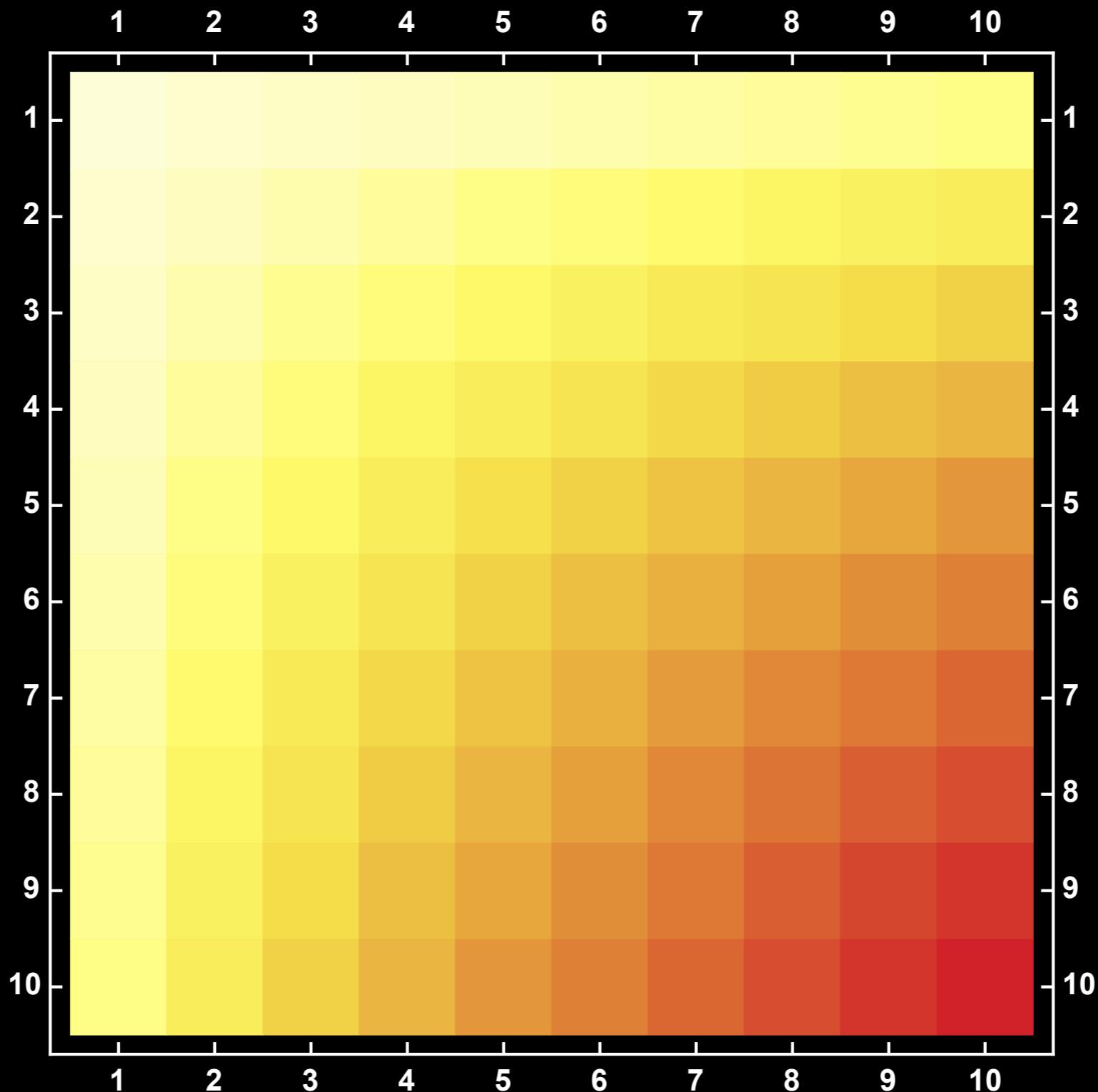


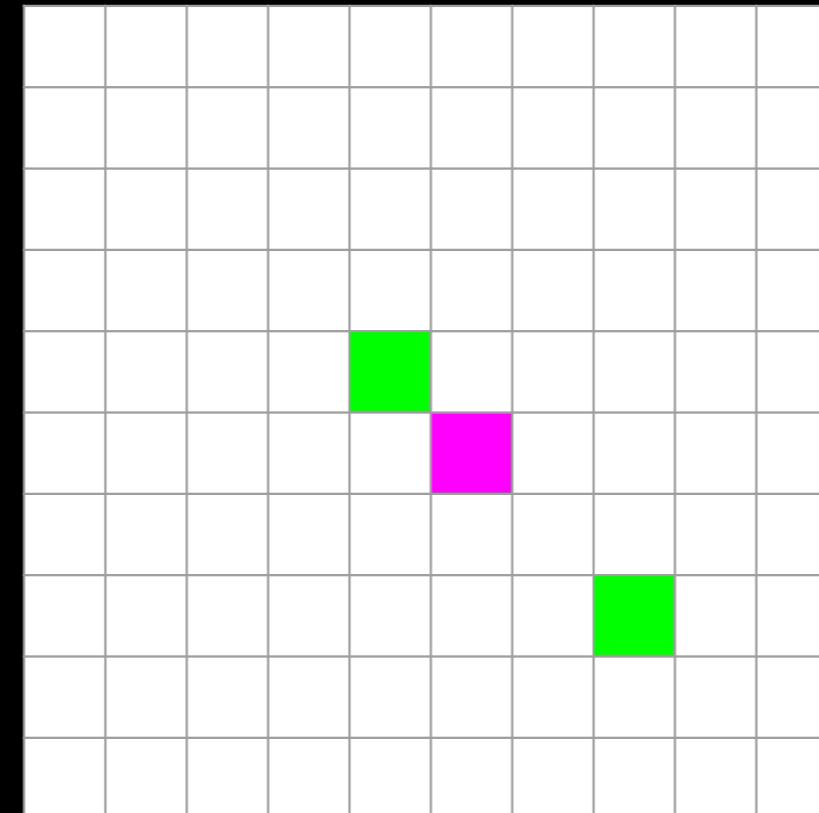
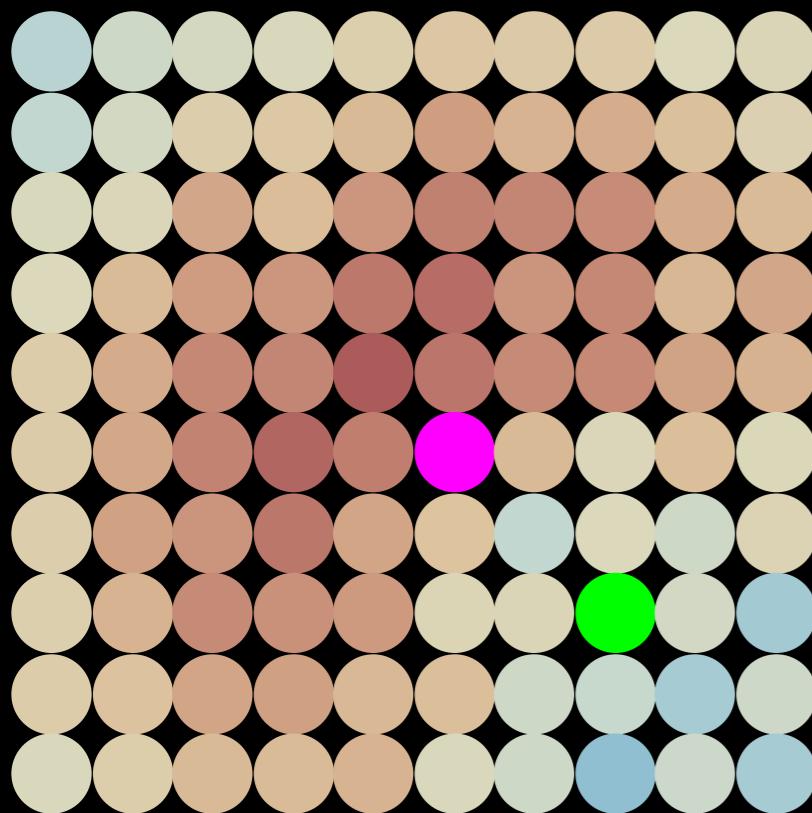
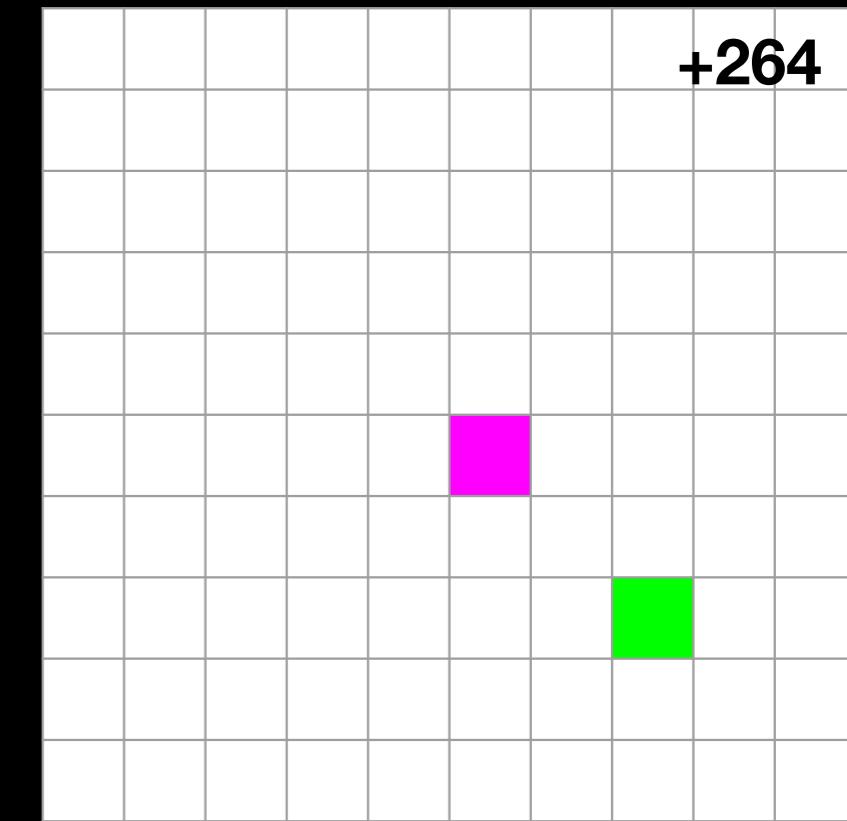
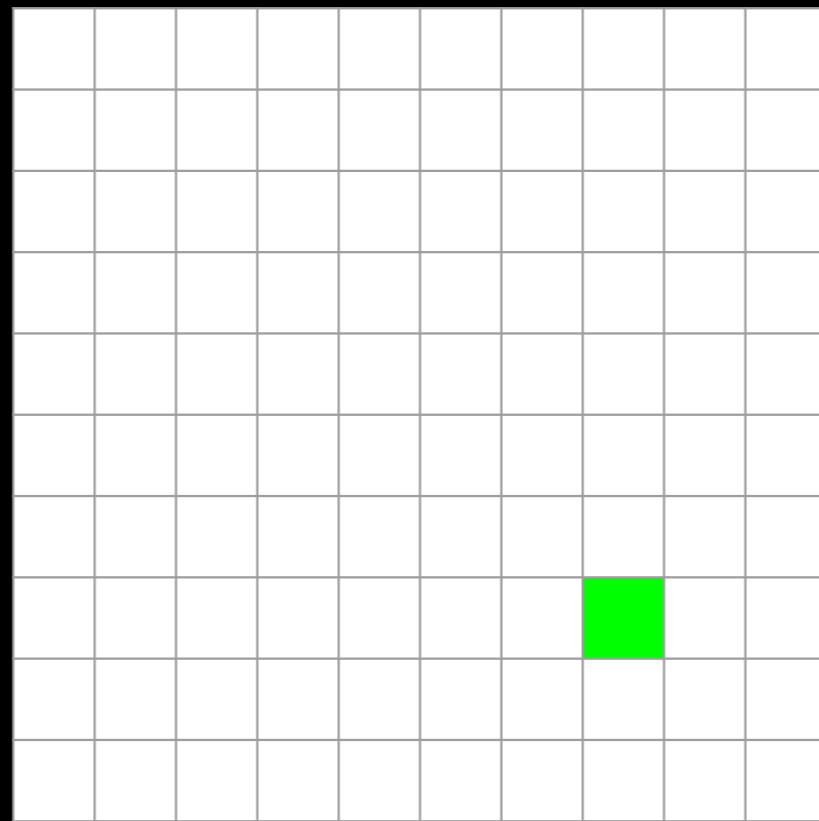
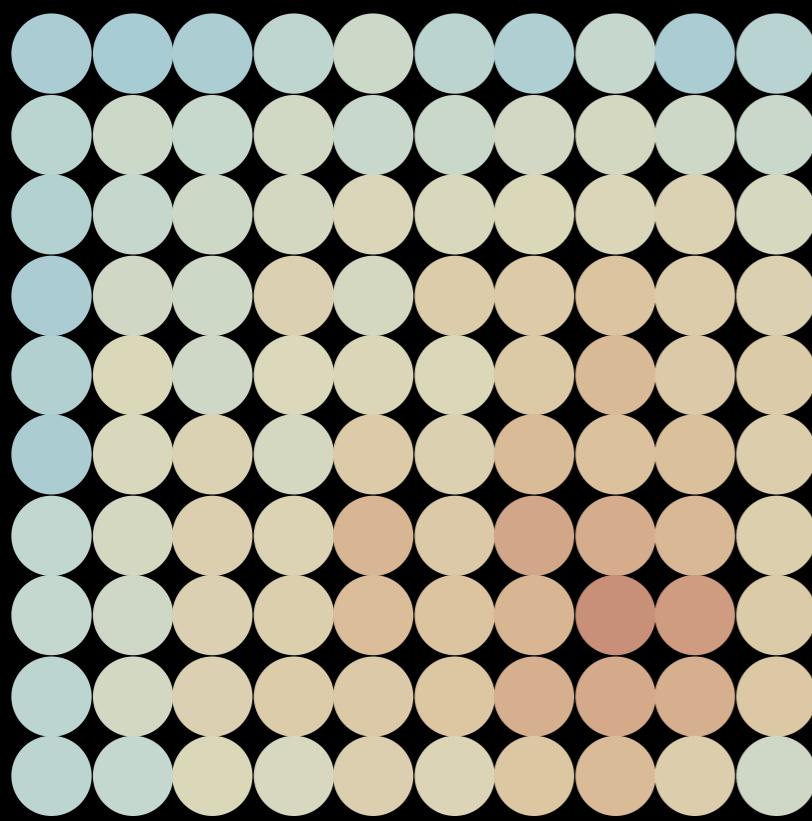
-9

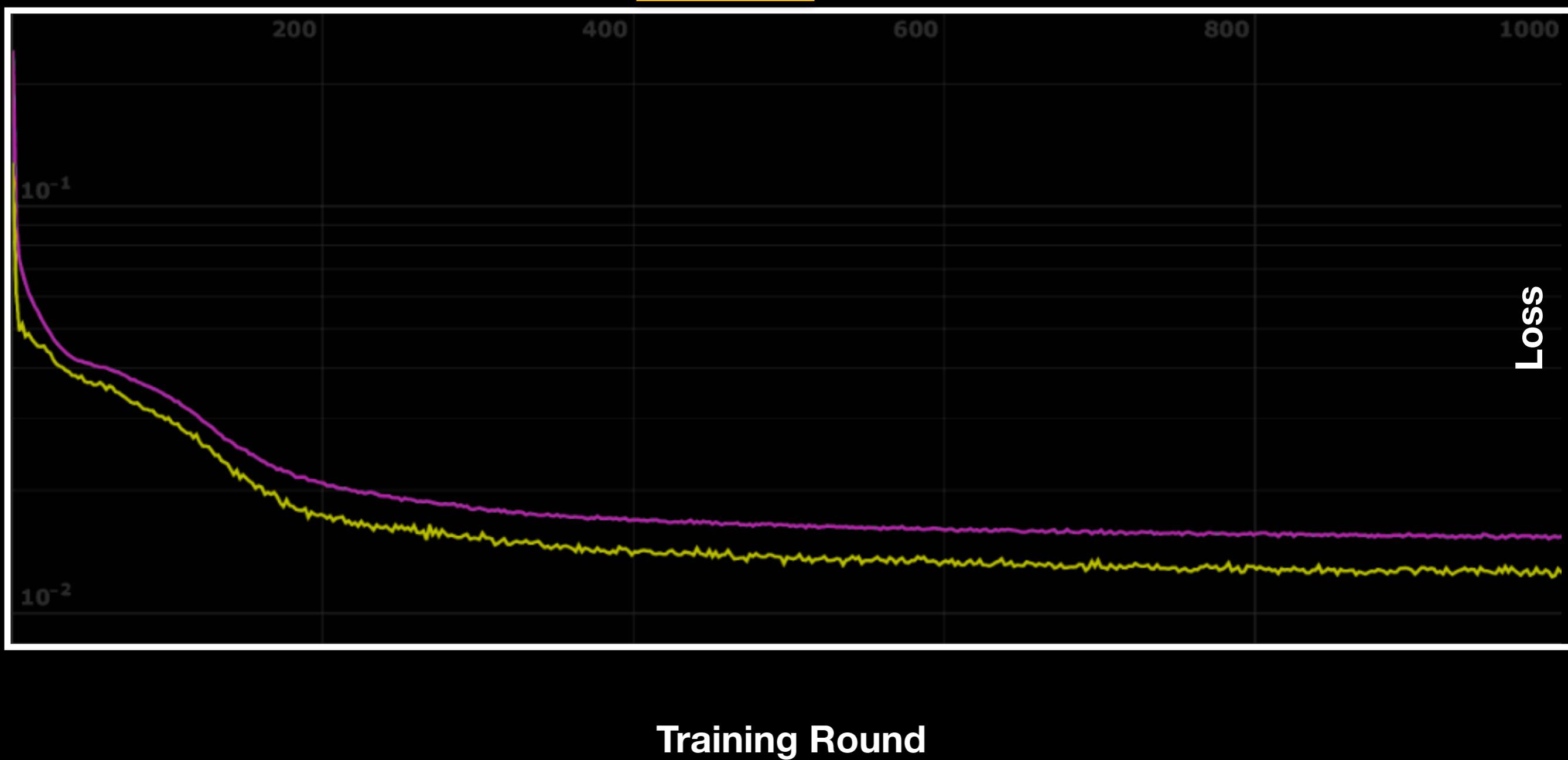
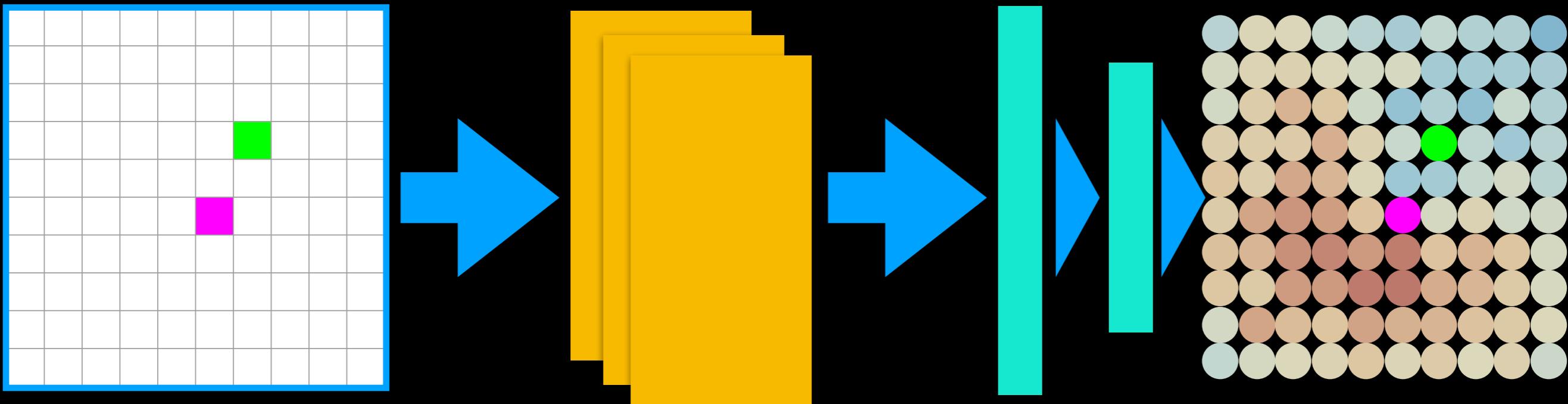


-2

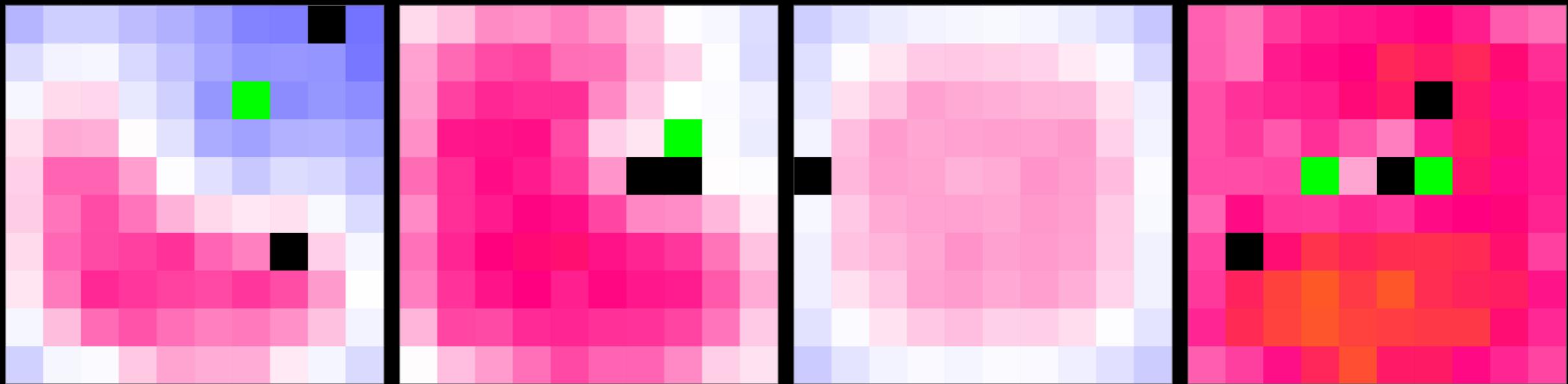
DEMO







Predicted



Target

