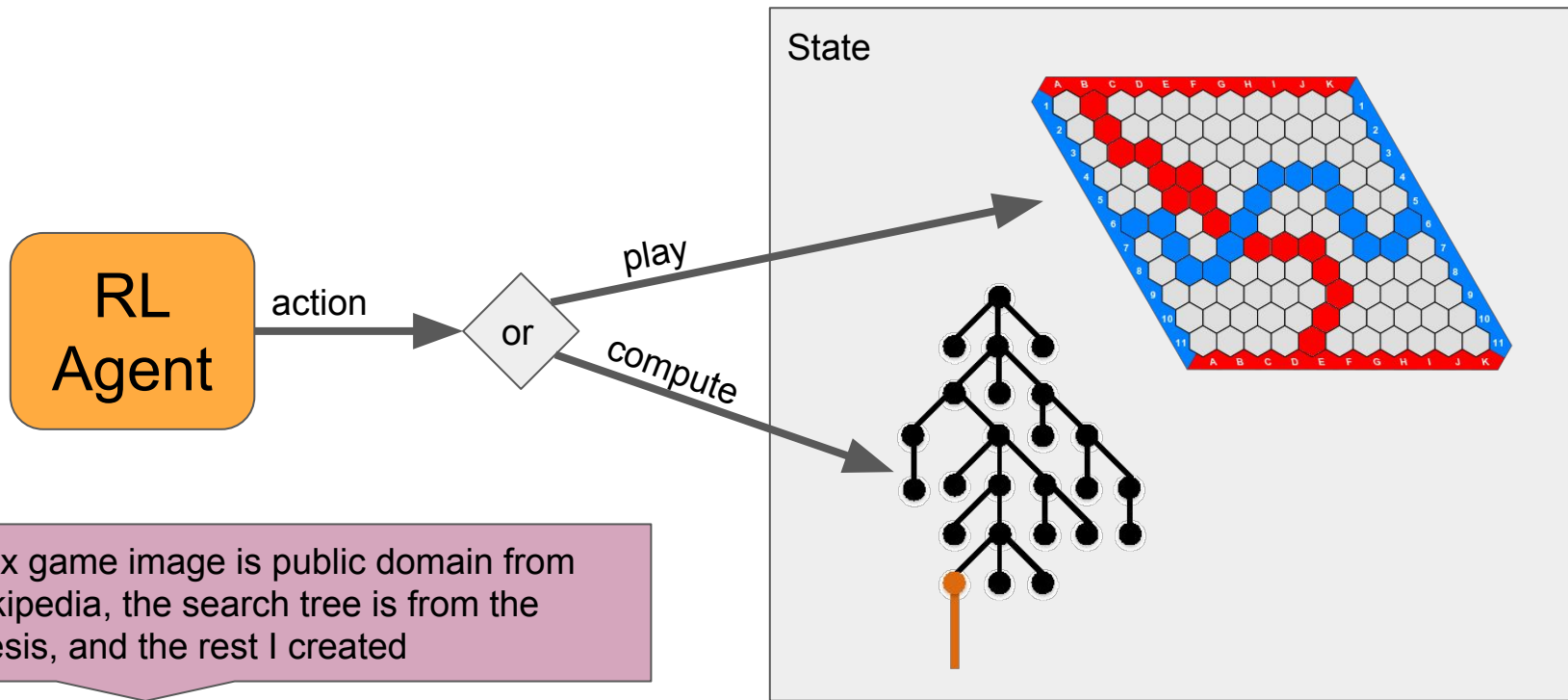


# Metalevel RL for MCTS

*Principles of Metalevel Control*  
*by Nicholas Hay*

Presented by Eric Langlois

# Metalevel Monte Carlo Tree Search



Hex game image is public domain from wikipedia, the search tree is from the thesis, and the rest I created

# Monte Carlo Tree Search

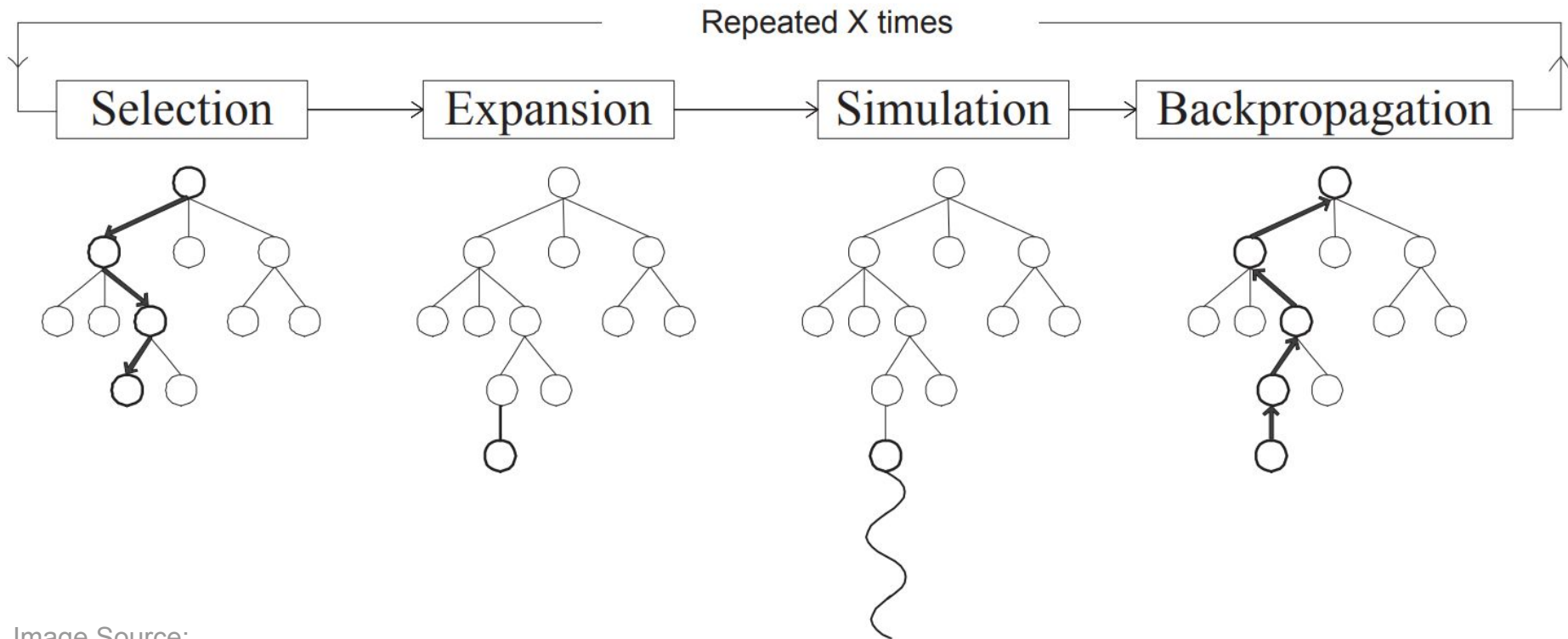
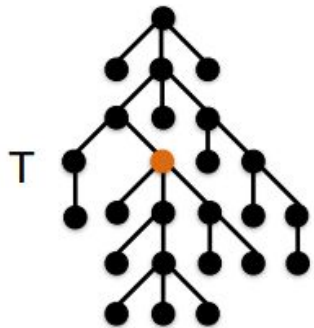


Image Source:

Chaslot et al (2008). Progressive Strategies for Monte-Carlo Tree Search. New Mathematics and Natural Computation.

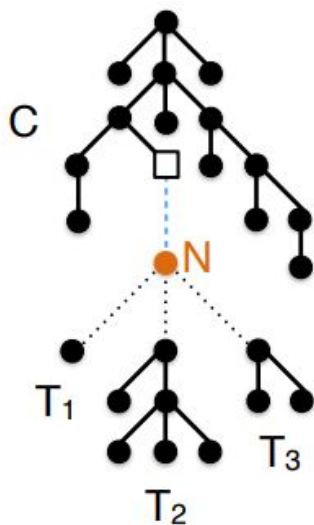
# Pointed Trees

From the thesis;  
I repositioned the labels



$T$

$$T = \text{PointedTree}(C, N, T_1, T_2, T_3).$$



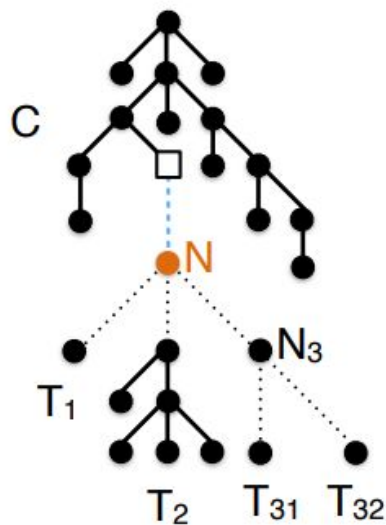
$C$

$T_1$

$T_2$

$T_3$

$$T_3 = \text{Tree}(N_3, T_{31}, T_{32}).$$



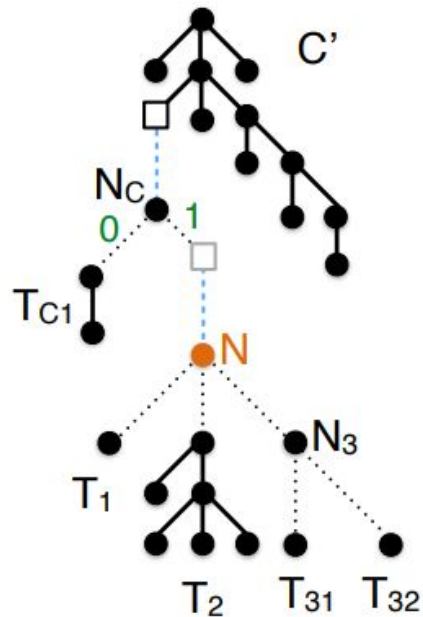
$C$

$T_1$

$T_2$

$T_{31}$

$T_{32}$



$C'$

$N_c$

$T_1$

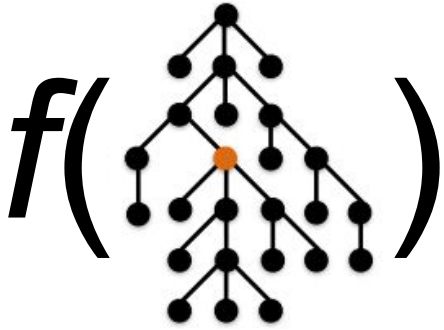
$T_2$

$T_{31}$

$T_{32}$

$$C = \text{Context}(C', N_c, 1, T_{c1})$$

# Recursive Functions on Pointed Trees



=

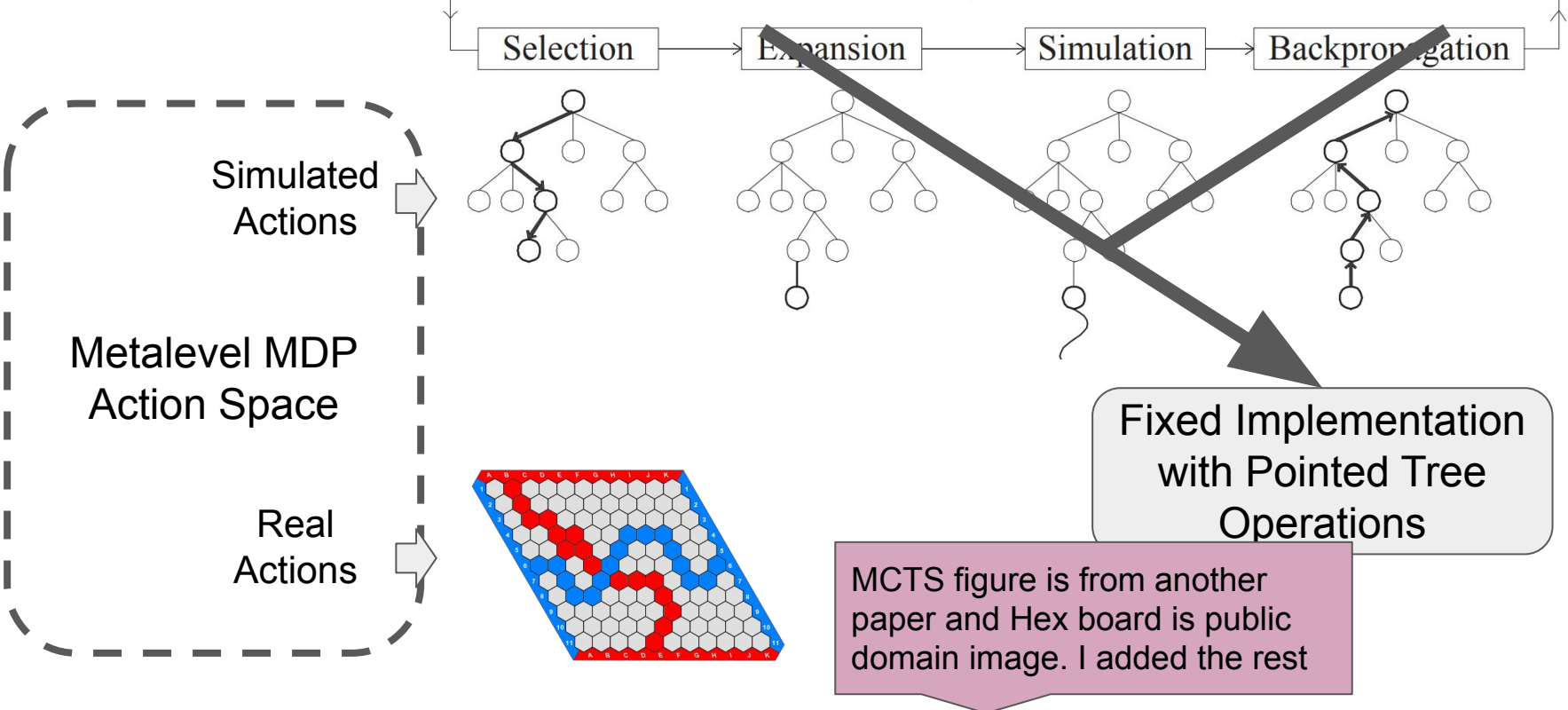
I created the formula display using copied bits of the tree image from the thesis

$(f_\theta, \frac{\partial f_\theta}{\partial \theta})$  is also recursive!

$$f_T(\text{Pointer Node}, f(\text{Context}), f(\text{1st Subtree}), f(\text{2nd Subtree}), \dots)$$



# MCTS as a Metalevel MDP



# Learning a Metalevel Agent

$$\pi_{\theta}(T) = g_{\theta}(f_{\text{fixed}}(T))$$

## *Tree Functions*

- num\_visits(T)
- all\_done(T)
- average\_rollout\_value(T)
- average\_estimate\_value(T)
- p\_over\_n+1(T)
- minimax(T)

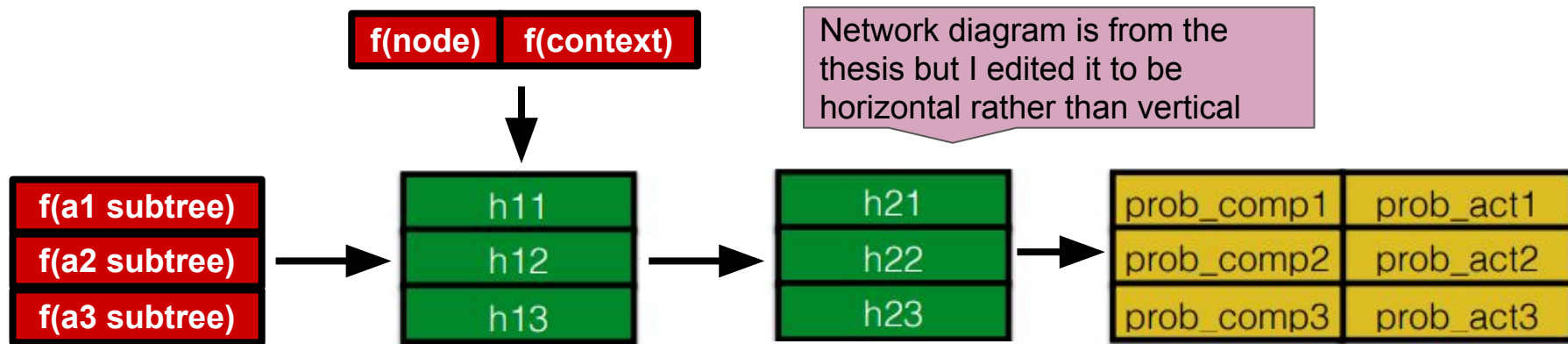
## *Context Functions*

- depth(C)
- alpha(C) *utility lower bound*
- beta(C) *utility upper bound*



# Policy Network $g_\theta$

$$\pi_\theta(T) = g_\theta(f_{\text{fixed}}(T))$$



## Training Algorithm: TRPO with Generalized Advantage Estimation

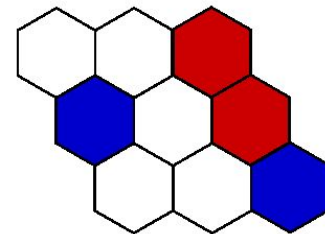
John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel.

Highdimensional continuous control using generalized advantage estimation. ICLR, 2016

# Experiments

Hex with a  $L \times L$  board, maximum of  $n$  simulated actions per move

Average win rate against a UCT baseline



Random Initialization

	n=10	n=20	n=50	n=100
L=3	<b>0.60</b>	<b>0.44</b>	-0.59	-0.60
L=5	<b>0.50</b>	<b>0.48</b>	-0.42	-0.91
L=7	<b>0.35</b>	<b>0.50</b>	-0.63	-0.54

Initialized to UCT

	n=10	n=20	n=50	n=100
L=3	<b>0.58</b>	<b>0.47</b>	<b>0.47</b>	<b>0.33</b>
L=5	<b>0.70</b>	<b>0.60</b>	<b>0.51</b>	-0.89
L=7	<b>0.37</b>	<b>0.47</b>	<b>0.36</b>	-0.45

# Related Work

## Metareasoning and Bounded Optimality

- Stuart J. Russell and Eric H. Wefald. Decision-theoretic control of search: General theory and an application to game-playing.
- Stuart Russell. *Rationality and intelligence: A brief update*. In Fundamental Issues of Artificial Intelligence, pages 7–28. Springer, 2014
- Eric Horvitz. *Models of continual computation*. In AAAI/IAAI, pages 286–293, 1997.

## Monte-Carlo Tree Search

- Levente Kocsis and Csaba Szepesvari. *Bandit Based Monte-Carlo Planning*. ECML, 2006.