



上海人工智能实验室  
Shanghai Artificial Intelligence Laboratory

# CBTune: Contextual Bandit Tuning for Logic Synthesis

Fangzhou Liu<sup>1</sup>, Zehua Pei<sup>1</sup>, Ziyang Yu<sup>1</sup>, Haisheng Zheng<sup>2</sup>,  
Zhuolun He<sup>1,2</sup>, Tinghuan Chen<sup>3</sup>, Bei Yu<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong, Hong Kong, China, <sup>2</sup>Shanghai Artificial Intelligence Laboratory, Shanghai, China <sup>3</sup>The Chinese University of Hong Kong, Shenzhen, China



## Highlights

- We propose **CBTune**, adapting the contextual bandit algorithm to facilitate efficient transformation selection through iterative model tuning.
- We implement the Syn-LinUCB algorithm as the agent and establish a context generator for informed decision-making in the bandit model.
- We present a novel “return-back” mechanism that revisits decisions to avoid local optima, distinguishing it from typical RL scenarios.
- Our method surpasses SOTA approaches for metrics and runtime within the same action space.

## Background

**ML-Enhanced Synthesis Optimization** Machine learning facilitates technology-independent optimization: 1) it models circuit structures to accurately predict performance metrics [4], 2) it employs reinforcement learning for rapid synthesis flow generation in an exponentially large solution space [1].

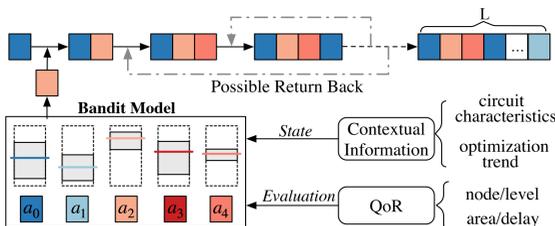


Figure 1. Illustration of our proposed contextual bandit-based approach for efficient synthesis flow generation.

**Bandit-based Search Model** The Multi-Arm Bandit (MAB) model, known for its efficiency in generating synthesis flows [3], strikes a balance between exploration and exploitation to optimize rewards. CBTune leverages domain-specific knowledge by integrating contextual data into the MAB model, enabling progressive decision-making depicted in Figure 1.

## Motivation

### Existing Problems

- NN-based methods are limited by time-consuming dataset preparation and training, as well as restricted transferability and system integration.
- The non-contextual MAB approach neglects key arm features like optimization trends and AIG characteristics. It also makes sequence-based decisions without considering permutations, compromising final performance.

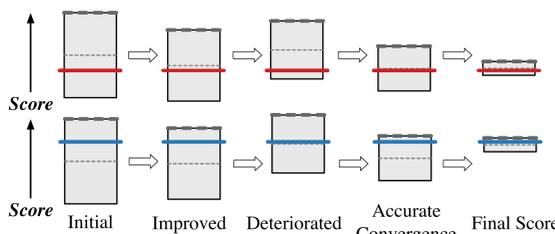


Figure 2. Score iterations for each arm in bandit model.

**Observations** LinUCB [2] improves MAB model by integrating contextual details like arm and environmental features to guide decision-making. The score for each arm  $a$  is updated by:

$$\text{LinUCB}_a = E(a|\mathbf{x}) + \alpha \text{STD}(a|\mathbf{x}) \\ = \mathbf{x}^\top \cdot \boldsymbol{\theta}_a + \alpha \sqrt{\mathbf{x}^\top \mathbf{A}_a^{-1} \mathbf{x}}. \quad (1)$$

### 1st term: Estimated Payoff

- Estimates average payoff from  $\mathbf{x}$
- $\boldsymbol{\theta}_a$  represents historical success

### 2nd term: Upper Confidence Bound

- Controlled by hyperparameter  $\alpha$
- Reflects uncertainty in estimation

Therefore, we propose a tailored bandit model to guide decisions for each individual transformation within the synthesis flow efficiently. This model:

1. Treats each transformation as an “arm” with equal initial UCB scores.
2. Iteratively updates scores to gauge performance.
3. Chooses and refines the highest-scoring arm in each iteration for enhanced score accuracy and reliability.
4. Steers scores towards the arms’ true payoffs, with the highest-scoring arm reflecting the best optimization performance.

## Pipeline

The overall CBTune framework is shown in Fig. 3.

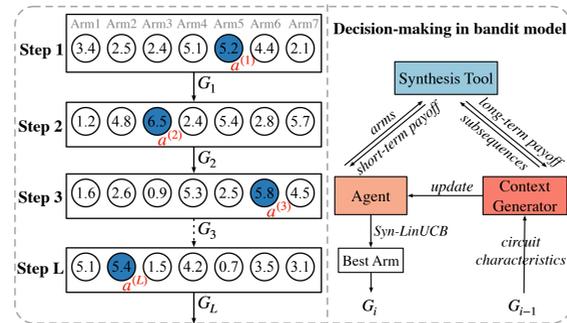


Figure 3. CBTune framework overview.

- **Action Space:**  $\mathcal{A} = \{\text{resub } (rs), \text{resub } -z \text{ (rsz)}, \text{rewrite } (rw), \text{rewrite } -z \text{ (rwz)}, \text{refactor } (rf), \text{refactor } -z \text{ (rfz)}, \text{balance } (b)\}$ .
- **Reward  $r$ :** the scaled payoff of a single arm execution.

## Methodology

**Context Generator** The vector  $\mathbf{x}$ , fusing circuit characteristics  $\mathbf{x}^c$  and the arm’s long-term payoff  $\mathbf{x}^l$ , informs the agent’s decisions by providing essential environmental and state insights.

Table 1. Contextual Information.

Feature	Example
Circuit Characteristics	Extracted by <code>yosys</code> and <code>ccirc</code> #Number of wires/cells/notes, #Maximum delay, #Number of combinational nodes, #Number of high degree comb, #Reconvergence, #Node shape...
Long-term Payoff of the Arm ( $\mathbf{x}^l$ )	Arm: <code>rewrite (rw)</code> ; $l = 5$ ; $m = 1$ ; [ <code>rw,rf,rf,rw,b</code> ] $\rightarrow$ Nodes: 28010, Level: 66 Arm: <code>refactor (rf)</code> ; $l = 4$ ; $m = 2$ ; [ <code>rf,b,rf,rw</code> ] $\rightarrow$ Nodes: 28350, Level: 69 [ <code>rf,rw,b,rs</code> ] $\rightarrow$ Nodes: 28324, Level: 67

**Agent: Syn-LinUCB** Key advantages:

1. It utilizes short-term payoffs to direct the agent to select arms toward the optimal target value per step, enhancing local performance.
2. It accounts for long-term payoffs to avert local optima and explore potential optimization trends, fostering improved decision quality.

### Algorithm 1 Syn-LinUCB

**Input:** Arms  $a \in \mathcal{A}$ , Context weights  $\mathbf{w} \in \mathbb{R}^d$ , Number of iterations  $T$ , Constant  $\rho$ .

**Output:** Best arm  $a_{best}$  in this step.

- 1:  $r_a \leftarrow$  Reward of all arms;
- 2: Extract the AIG characteristics:  $\mathbf{x}_a^c \in \mathbb{R}^{d_1}$ ;
- 3: Arm selection times  $s_a = 0$ ;
- 4: **for**  $t = 1, 2, \dots, T$  **do**
- 5: Update the long-term payoff:  $\mathbf{x}_{t,a}^l \in \mathbb{R}^{d_2}$ ;
- 6: Observe features of  $a \in \mathcal{A}$ :  $\mathbf{x}_{t,a} = [\mathbf{x}_{t,a}^c, \mathbf{x}_{t,a}^l] \in \mathbb{R}^d$ ;
- 7: **for**  $\forall a \in \mathcal{A}$  **do**
- 8: Initialize historical context and reward by  $\mathbf{A}_a = \mathbf{I}_d$ ,  $\mathbf{b}_a = \mathbf{0}_d$ ,  $\forall a$  is new;
- 9: Update hyperparameter  $\alpha$  by  $\alpha = 1.0 + \frac{\log(2.0/\rho)}{s_a}$ ;
- 10: Update the decision parameter by  $\boldsymbol{\theta}_a = \mathbf{A}_a^{-1} \mathbf{b}_a$ ;
- 11: Calculate the weighted context  $\mathbf{x}_{t,a}^w = \mathbf{x}_{t,a} \mathbf{w}$ ;
- 12: Update score by  $p_{t,a} = \boldsymbol{\theta}_a^\top (\mathbf{x}_{t,a}^w) + \alpha \sqrt{(\mathbf{x}_{t,a}^w)^\top \mathbf{A}_a^{-1} (\mathbf{x}_{t,a}^w)}$ ;
- 13: **end for**
- 14: Choose arm by  $a_t = \arg \max_{a \in \mathcal{A}} p_{t,a}$ ;
- 15: Increase the selection count of arm  $a_t$  by  $s_{a_t} = s_{a_t} + 1$ ;
- 16: Update the parameters  $\mathbf{A}_{a_t}$  and  $\mathbf{b}_{a_t}$  of the chosen arm  $a_t$  by
- 17:  $\mathbf{A}_{a_t} = \mathbf{A}_{a_t} + \mathbf{x}_{t,a_t} \mathbf{x}_{t,a_t}^\top$ ,  $\mathbf{b}_{a_t} = \mathbf{b}_{a_t} + r_{a_t} \mathbf{x}_{t,a_t}$ ;
- 18: **end for**
- 19:  $a_{best} \leftarrow a_t$ .

**Return-back Mechanism** To amend suboptimal decisions stemming from a lack of historical data, we allow CBTune the capacity to “regret” by recording synthesis results in a hash table. This allows CBTune to compare new results with past decisions and, if necessary, return to a crucial step to reselect a better arm, thus improving decision quality.

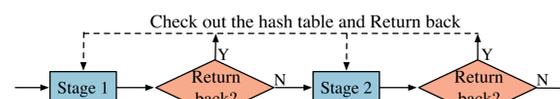


Figure 4. The return-back mechanism in CBTune.

## Evaluation Results

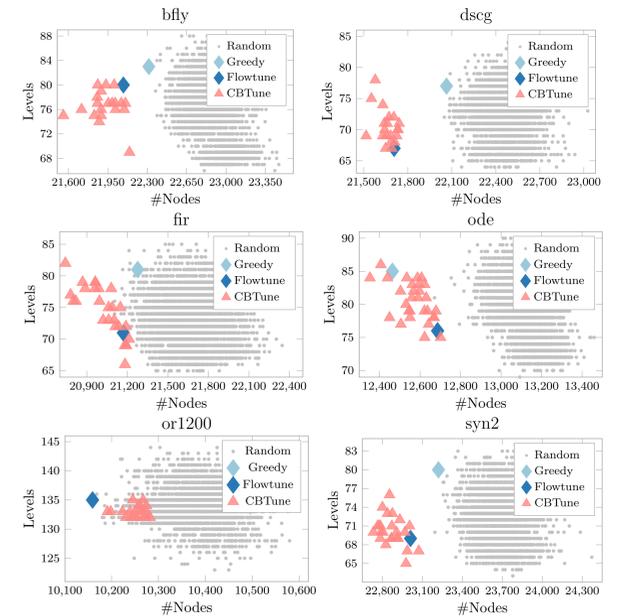


Figure 5. CBTune vs. FlowTune [3] in AIG node optimization.

Table 2. CBTune vs. FlowTune in 6-LUTs optimization.

Benchmark	Initial	Greedy	Flowtune [3]	CBTune			
	#LUTs	#LUTs	#LUTs	#LUTs	#LUTs	$\tau(m)$	
<i>bfly</i>	9019	8269	8216	76.47	<b>7962</b>	8086.03	<b>29.63</b>
<i>dscg</i>	8534	8313	8302	77.15	<b>7981</b>	8119.84	<b>30.44</b>
<i>fir</i>	8646	8385	8094	74.23	<b>7820</b>	7977.38	<b>27.6</b>
<i>ode</i>	5244	5316	5096	34.83	<b>4920</b>	5046.71	<b>17.32</b>
<i>or1200</i>	2776	2748	2747	20.08	<b>2731</b>	2754.07	<b>15.62</b>
<i>syn2</i>	8777	8669	8603	81.33	<b>8234</b>	8360.53	<b>31.67</b>
GEOMEAN	6631.20	6464.69	6364.89	54.04	<b>6166.39</b>	6271.82	<b>24.48</b>
Ratio Avg.	1.000	0.975	0.960	1.000	<b>0.930</b>	0.946	<b>0.453</b>

Table 3. CBTune vs. NN-enhanced RL in 6-LUTs optimization.

Benchmark	Initial	Greedy	DRILLS [1]		RL4LS*		CBTune	
	#LUTs	#LUTs	#LUTs	$\tau(m)$	#LUTs	$\tau(m)$	#LUTs	$\tau(m)$
<i>max</i>	721	697	694	32.58	687.8	54.34	<b>684.25</b>	<b>6.01</b>
<i>adder</i>	249	<b>244</b>	<b>244</b>	24.05	<b>244</b>	10.05	<b>244</b>	<b>5.97</b>
<i>cavlc</i>	116	115	112.2	26.02	111.3	3.22	<b>111</b>	<b>2.37</b>
<i>ctrl</i>	29	<b>28</b>	<b>28</b>	24.25	<b>28</b>	2.85	<b>28</b>	<b>0.59</b>
<i>int2float</i>	47	46	42.6	21.7	42.3	2.81	<b>40</b>	<b>2.76</b>
<i>router</i>	73	67	70.1	22.01	69.5	3.07	<b>68.11</b>	<b>2.32</b>
<i>priority</i>	264	146	<b>133.4</b>	23.32	142.9	5.9	138.86	<b>3.41</b>
<i>i2c</i>	353	291	292.1	25.17	289.32	7.55	<b>283.11</b>	<b>3.61</b>
<i>sin</i>	1444	1451	1441.5	51.15	<b>1438</b>	20.1	1441.67	<b>9.71</b>
<i>square</i>	3994	3898	3889.4	130	3889	72.88	<b>3882.11</b>	<b>25.99</b>
<i>sqrt</i>	8084	4807	4708	147.64	4685.3	196.15	<b>4607</b>	<b>36.51</b>
<i>log2</i>	7584	7660	7583.6	198.6	7580.1	125.28	<b>7580</b>	<b>41.27</b>
<i>multiplier</i>	5678	5688	5678	180.84	<b>5672</b>	187.81	5679.75	<b>29.08</b>
<i>voter</i>	2744	1904	1834.7	84.43	<b>1678.1</b>	330.48	1682.25	<b>11.46</b>
<i>div</i>	23864	4205	7944.4	259.75	7807.1	482	<b>4180.91</b>	<b>25.58</b>
<i>mem_ctrl</i>	11631	<b>10144</b>	10527.6	229.33	10309.7	1985.84	10242.57	<b>45.81</b>
GEOMEAN	926.59	732.69	753.49	59.48	748.34	34.54	<b>712.83</b>	<b>8.37</b>
Ratio Avg.	1.000	0.791	0.813	1.000	0.808	0.581	<b>0.769</b>	<b>0.141</b>

\* Last 10 in RL-PPO-Pruned [5]

## Conclusion

- CBTune outperforms FlowTune in both AIG nodes/6-LUT optimization in both metric and runtime. Our method also outshines three RL-based methods by reducing 6-LUT counts up to 4.4%, all achieved in a swift 8.37 minutes per design.
- CBTune efficiently generates synthesis flows with excellent, stable results and fast runtime, without training data or complex procedures.

## References

- [1] Abdelrahman Hosny, Soheil Hashemi, et al. DRILLS: Deep reinforcement learning for logic synthesis. pages 581–586, 2020.
- [2] Lihong Li, Wei Chu, et al. A contextual-bandit approach to personalized news article recommendation. pages 661–670, 2010.
- [3] Walter Lau Neto, Yingjie Li, et al. Flowtune: End-to-end automatic logic optimization exploration via domain-specific multi-armed bandit. 2022.
- [4] Nan Wu et al. Lostin: Logic optimization via spatio-temporal information with hybrid graph models. pages 11–18, 2022.
- [5] Guanglei Zhou and Jason H Anderson. Area-Driven FPGA Logic Synthesis Using Reinforcement Learning. pages 159–165, 2023.