

PrefixRL: Optimization of Parallel Prefix Circuits using Deep Reinforcement Learning

Rajarshi Roy, Jonathan Raiman, Neel Kant, Ilyas Elkin, Robert Kirby,
Michael Siu, Stuart Oberman, Saad Godil, Bryan Catanzaro
NVIDIA, Santa Clara, CA, USA

{rajarshir, jraiman, nkant, ielkin, rkirby, msiu, soberman, sgodil, bcatanzaro}@nvidia.com

Abstract—In this work, we present a reinforcement learning (RL) based approach to designing parallel prefix circuits such as adders or priority encoders that are fundamental to high-performance digital design. Unlike prior methods, our approach designs solutions tabula rasa purely through learning with synthesis in the loop. We design a grid-based state-action representation and an RL environment for constructing legal prefix circuits. Deep Convolutional RL agents trained on this environment produce prefix adder circuits that Pareto-dominate existing baselines with up to 16.0% and 30.2% lower area for the same delay in the 32b and 64b settings respectively. We observe that agents trained with open-source synthesis tools and cell library can design adder circuits that achieve lower area and delay than commercial tool adders in an industrial cell library.

Index Terms—machine learning, reinforcement learning, datapath optimization

I. INTRODUCTION

Several fundamental digital design building blocks such as adders, priority encoders, inc(dec)rementers and gray-to-binary code converters can be reduced to prefix-sum computations and implemented as (parallel) prefix circuits [1]. Thus, the optimization of prefix circuits for area, delay and power is an important and well studied problem in digital hardware design.

The optimization of prefix circuits is challenging as their large design space grows exponentially with input length and is intractable to enumerate. As a result, exhaustive search approaches do not scale beyond small input lengths [2]. Several regular prefix circuit structures [3]–[5] have been proposed that trade off logic level, maximum fanout and wiring tracks. Another set of algorithms [6]–[9] optimize prefix circuit size and level properties. However, [10] observes that prefix circuit level and maximum fanout properties do not map to circuit area, power and delay due to physical design complexities such as capacitive loading and congestion.

Meanwhile, there are a growing number of success stories in other domains using deep reinforcement learning (RL) [11] to produce sophisticated solutions to sequential decision problems. RL agents have outperformed humans in complex games [12] and produced novel solutions to search and design problems [13]. PrefixRL continues this trend in the circuit design domain. The key contributions of this work are:

- An RL framework that trains an agent, tabula rasa, to explore the unrestricted prefix circuit space with synthesis in the loop (Fig. 1) while optimizing for area and delay.
- We demonstrate the capabilities of PrefixRL with the task of optimizing 32b and 64b prefix adder circuits. PrefixRL explores the massive $\mathcal{O}(2^{N^2})$ design space and generates state-of-the-art frontiers of designs that Pareto-dominate all the designs found by prior work using simulated annealing [14], exhaustive search with pruning [15] [10] and regular structures [3]–[5] when synthesized in an industrial process technology, achieving

a maximum area savings of 16.0% and 30.2% for equivalent delay targets in the 32b and 64b settings.

- We observe that agents trained with open-source synthesis tools (OpenPhySyn [16]) and cell library (Nangate45 [17]) design adder circuits that can achieve lower area and delay than commercial tool adders in an industrial cell library (8nm).
- We illustrate the challenges of training an RL agent with continuous synthesis feedback and the solutions required to handle workloads in similar use cases. With our state-of-the-art results, these solutions provide a compelling blueprint for future RL approaches to design automation.

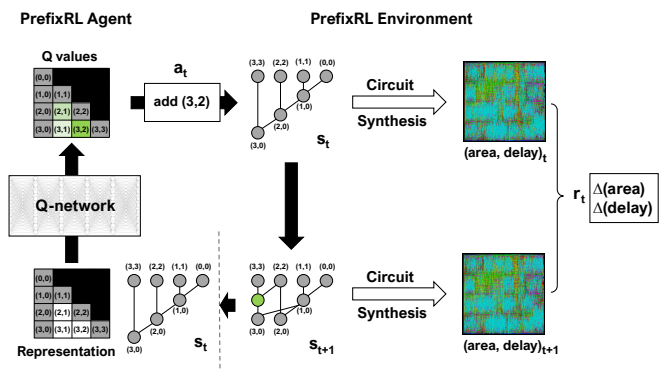


Fig. 1: PrefixRL flow. s_t shown is the ripple-carry prefix graph, a possible starting state s_0 . The Q-network takes action (3,2) modifying the circuit and receiving a reward computed by the area/delay difference in the circuits corresponding to s_t and s_{t+1} . Refer to Section IV-D for circuit synthesis details. For clarity we show a 4b circuit but our primary results are obtained in the 32b and 64b setting.

II. RELATED WORK

Several approaches attempt to optimize the prefix circuits directly for area, power and delay. [15] utilizes a combination of heuristic rules to prune the intractable design space of 8, 16, 32 and 64 bit adders to a small subset that can be exhaustively searched. [10] extends this technique by proposing an alternative set of pruning heuristics that result in a larger set of pruned adders which are then searched using a machine learning model that is trained to predict physical metrics. Our approach is fundamentally different as it does not rely on any hand-crafted heuristics. Our agents learn tabula rasa to explore the full prefix adder space with synthesis in the loop.

Another work [14] uses simulated annealing (SA) [18] to randomly modify and legalize 32b prefix adder graphs towards optimal structures in an unrestricted design space evaluated using an analytical model. They obtain a frontier of designs that optimize these competing metrics more effectively than the regular structures [3]–[5] as well as those obtained by the previously mentioned exhaustive search on a heuristically pruned search space [15]. However, all of

these results are limited to analytical evaluation metrics. Our prefix RL approach instead optimizes physical synthesis metrics directly with synthesis in the loop. We investigate applying our approach with the analytical evaluation metrics and demonstrate that our prefix circuits Pareto dominate the SA based baseline circuits (Section V-D). However, we observe that these circuits that were found using the analytical evaluation metric significantly degrade in quality once they go through physical synthesis. Note that since physical synthesis is significantly more computationally intensive than the analytical evaluation, it would not be feasible to scale a fundamentally sequential algorithm such as SA to use physical synthesis in the loop.

III. BACKGROUND

A. Prefix Graphs

The generalized prefix-sum computation is to compute $y_i = x_i \circ x_{i-1} \circ \dots \circ x_0$ for $0 \leq i \leq N-1$, given N inputs x_0, x_1, \dots, x_{N-1} and a binary associative operator \circ [19].

An N -input prefix-sum computation can be performed in several ways due to the associativity of the operator. For example, two of the ways the 4-input prefix sum can be computed are:

$$\begin{aligned} y_0 &= x_0, y_1 = x_1 \circ y_0, y_2 = x_2 \circ y_1, y_3 = x_3 \circ y_2 \\ y_0 &= x_0, y_1 = x_1 \circ y_0, y_2 = x_2 \circ y_1, z_{3:2} = x_3 \circ x_2, y_3 = z_{3:2} \circ y_1 \end{aligned}$$

Introducing the additional term $z_{3:2}$, breaks the dependency of y_3 on y_2 and allows it to be computed in parallel with y_2 , thus the term parallel prefix [20]. In general, we denote $z_{i:j}$ to represent $x_i \circ x_{i-1} \circ \dots \circ x_j$. Then the outputs y_i can be rewritten as $z_{i:0}$ and inputs x_i can be rewritten as $z_{i:i}$. Note that y_0 and x_0 both correspond to $z_{0:0}$.

Parallel prefix computations can be represented as a directed acyclic prefix graph where every computation unit $z_{i:j}$ is a graph node that performs a single operation on a pair of inputs: $z_{i:j} = z_{i:k} \circ z_{k-1:j}$ where $i \geq k > j$. We use the notation from [15] where the most and least significant bits (*MSB*, *LSB*) of computation node $z_{i:j}$ is (i, j) . Using this notation we will term the node (i, k) as the *upper parent* of (i, j) and the node $(k-1, j)$ as its *lower parent*. The prefix graphs corresponding to the 4-input prefix sum computations above are shown in Fig. 1 as s_t and s_{t+1} . In both graphs, the upper and lower parents of node $(2, 0)$ are $(2, 2)$ and $(1, 0)$.

Every legal N -input prefix graph must have input nodes (i, i) , output nodes $(i, 0)$ for $1 \leq i \leq N-1$, and the input/output node $(0, 0)$. Furthermore, every non-input node must have exactly one upper parent (*up*) and one lower parent (*lp*) such that:

$$\begin{aligned} LSB(node) &= LSB(lp(node)) \\ LSB(lp(node)) &\leq MSB(lp(node)) \\ MSB(lp(node)) &= LSB(up(node)) - 1 \\ LSB(up(node)) &\leq MSB(up(node)) \\ MSB(up(node)) &= MSB(node) \end{aligned} \quad (1)$$

B. Deep Reinforcement Learning

Reinforcement learning (RL) [11] is a class of algorithms applicable to sequential decision making tasks. RL makes use of the Markov Decision Process (MDP) formalism wherein an *agent* attempts to optimize a function in its *environment*. An MDP can be completely described by a *state space* \mathcal{S} (with states $s \in \mathcal{S}$), an *action space* \mathcal{A} (with actions $a \in \mathcal{A}$), a *transition function* $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ and a *reward function* $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. In an MDP, an episode evolves over discrete timesteps $t = 0, 1, 2, \dots$ where the agent observes s_t and responds with action a_t using a *policy* $\pi(a_t|s_t)$. The environment provides to the agent the next state $s_{t+1} = \mathcal{T}(s_t, a_t)$ and the reward

$r_t = \mathcal{R}(s_t, a_t)$. The agent is tasked with maximizing the *return* (cumulative future rewards) by learning an optimal policy π^* .

The Q value of a state-action pair (s_t, a_t) under a policy π is defined to be the expected return if action a_t is taken at state s_t and future actions are taken using the policy π .

$$Q^\pi(s_t, a_t) = \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots], \gamma \in [0, 1] \quad (2)$$

The discount factor $\gamma \in [0, 1]$ balances short-term versus long-term rewards. The Q -learning algorithm [21] starts the agent with a random policy and uses the experience gathered during its interaction with the environment (s_t, a_t, r_t, s_{t+1}) to iterate towards an optimal policy by updating Q with a learning rate $\alpha \in [0, 1]$:

$$Q(s_t, a_t) \leftarrow (1 - \alpha) * Q(s_t, a_t) + \alpha * (r_t + \gamma \max_{a'} Q(s_{t+1}, a')) \quad (3)$$

The policy for a Q -learning agent is simply $\pi(\cdot|s_t) = \text{argmax}_a Q(s_t, a)$. We use the ϵ -greedy policy, where random actions are chosen with a probability ϵ to increase exploration in the state space. ϵ is annealed to zero during the course of training and is always zero when doing evaluation.

The DQN (deep Q network) algorithm [22] uses a deep neural network as a Q value function approximator to achieve human-level performance on several Atari games. DQN stabilizes training using a second *target network* to estimate the Q values of (s_{t+1}, a') that is updated less frequently and sampling an *experience replay* buffer. The Double-DQN algorithm [23] further improves training by reducing harmful overestimations in DQN.

IV. PREFIXRL IMPLEMENTATION

We frame the optimization of prefix circuits as a RL task by creating an MDP for their construction. We pick prefix adders due to their importance in arithmetic datapaths and focus on minimizing circuit area and delay. We then train multiple PrefixRL agents to design an area-delay minimized Pareto frontier of adders.

A. Reinforcement Learning Environment

The PrefixRL state space \mathcal{S} consists of all legal N -input prefix graphs. N -input graphs can be represented in a $N \times N$ grid with rows representing *MSB* and columns representing *LSB* (Fig. 1). Note that the input nodes ($MSB = LSB$) will lie on the diagonal, output nodes will lie on the first column ($LSB = 0$) and locations above the diagonal ($LSB > MSB$) cannot contain a node. The remaining $(N-1)(N-2)/2$ locations where non-input/output nodes may or may not exist define the $\mathcal{O}(2^{(N-1)(N-2)/2}) = \mathcal{O}(2^{N^2})$ state space of N -input prefix graphs. For example, 32-input graphs will have $|\mathcal{S}| = \mathcal{O}(2^{465})$ with a lower exact value because not all combinations of nodes in those locations will meet the legality constraints in (1).

The action space \mathcal{A} for an N -input prefix graph consists of two actions (add or delete) for any non-input/output node i.e. where $LSB \in [1, N-2]$ and $MSB \in [LSB+1, N-1]$. Hence, $|\mathcal{A}| = (N-1)(N-2)/2$. The environment evolution through \mathcal{T} always maintains a legal prefix graph by:

- 1) Applying a legalization procedure after an action that may add or delete additional nodes to maintain legality.
- 2) Forbidding redundant actions that gets undone by the legalization procedure.

During legalization, the upper parent of a node, $up(node)$, is the existing node with same *MSB* and the next highest *LSB*. The lower parent of a node is computed using the node and its upper parent (1):

$$(MSB_{lp(node)}, LSB_{lp(node)}) = (LSB_{up(node)} - 1, LSB_{node})$$

An illegal condition happens only when the lower parent $lp(node)$ of a node does not exist, so the legalization procedure adds any missing lower parent nodes.

Algorithm 1: PrefixRL N-input prefix graph actions

```

Function Initialize:
  nodelist ← ∅, minlist ← ∅;
  for m ← 0 to (N - 1) do
    | add (m, m), (m, 0) to nodelist // add in/out nodes
  end
Function Add (msb,lsb):
  add (msb, lsb) to minlist;
  // remove new node's and child's lps from minlist
  for l ← (msb - 1) to 0 do
    | if (msb, l) is in minlist then
      | | delete lp(msb, l) from minlist
    end
  end
  Legalize()
Function Delete (msb,lsb):
  delete (msb, lsb) from minlist;
  Legalize()
Function Legalize:
  nodelist ← minlist;
  for m ← 0 to (N - 1) do
    | add (m, m), (m, 0) to nodelist // add in/out nodes
  end
  for m ← (N - 1) to 0 do
    | for l ← (m - 1) to 0 do
      | | if (m, l) is in nodelist then
        | | | add lp(m, l) to nodelist
      end
    end
  end
end

```

The action of adding a node that already exists (in *nodelist*) is redundant and is forbidden. Deleting is limited to nodes in *minlist* (nodes that are not lower parents of other nodes) to prevent legalization from adding back deleted nodes (Algorithm 1).

B. Scalarized Double Deep Q Learning

PrefixRL uses a scalarized version of the Double-DQN algorithm [23]. Every episode starts the environment with the initial state s_0 randomly chosen to be either a ripple-carry or a Sklansky [3] graph. These are prefix graphs with the minimum node and level count respectively. Every action a_t from the agent modifies the legal prefix graph s_t to another legal prefix graph s_{t+1} and returns a reward vector \mathbf{r}_t that indicates the decrease in the normalized circuit area and delay when its prefix graph is modified from s_t to s_{t+1} (Fig. 1). Details of how we measure area and delay are given in Section IV-D.

$$\mathbf{r}_t = [area(s_t) - area(s_{t+1}), delay(s_t) - delay(s_{t+1})]$$

With competing objectives such as area and delay, the same improvement in the scalarized objective can occur from either an improvement in area or delay, but the resulting prefix graph structures would be very different. Thus, it is difficult for a Q-learning agent to infer how its actions affect prefix graphs if it can only observe a scalarized objective. The *scalarized deep Q-learning* algorithm [24] updates the Q-learning algorithm for such multi-objective settings by receiving rewards and learning the Q value for different objectives separately but choosing actions after scalarizing the Q values with a weight vector \mathbf{w} .

The PrefixRL agent estimates the Q function vector $\mathbf{Q}(s, a) = [Q_{area}(s, a), Q_{delay}(s, a)]$ for a pair of prefix graph state and modification action using a deep neural network. The double-DQN training procedure [23] is used with the target, loss and action selection extended for scalarization:

$$\mathbf{y}_t = \mathbf{r}_t + \gamma \mathbf{Q}(s_{t+1}, \arg\max_a [\mathbf{w}^\top \mathbf{Q}(s_{t+1}, a; \theta_t)]; \theta'_t) \quad (4)$$

$$Loss(\mathbf{Q}(s_t, a_t; \theta_t), \mathbf{y}_t) \quad (5)$$

$$a_t = \arg\max_a [\mathbf{w}^\top \mathbf{Q}(s_t, a; \theta_t)] \quad (6)$$

Where θ_t and θ'_t are the parameters of the local and target networks in the double-DQN algorithm.

A Pareto frontier of designs can then be obtained by solving multiple single-objective optimization problems, each with the scalarized objective $\mathbf{w}^\top \mathbf{m} = w_{area} \cdot area + w_{delay} \cdot delay$ using different scalarization weights \mathbf{w} . We conceptually encode a tradeoff of objectives by normalizing \mathbf{w} such that its elements are nonnegative and sum to 1. However, we must also scale the raw values of *area* and *delay* since their units are incomparable. Our procedure for this is to multiply those values by scaling constants c_{area}, c_{delay} such that the Pareto frontier for different \mathbf{w} evenly covers the breadth of baseline prefix graph designs synthesized with multiple delay targets. In our experiments we use $c_{area} = 0.001$ and $c_{delay} = 10$.

C. Q Network Architecture

The deep neural network Q value approximator takes the state s_t as the input and predicts $\forall a \in \mathcal{A} : [Q_{area}(s_t, a), Q_{delay}(s_t, a)]$. Based on the $N \times N$ grid based representation of prefix graphs described in Section IV-A, the input to the neural network is a $N \times N \times 4$ tensor where the 4 channels encode node features as:

- 1) 1 if node (MSB, LSB) in *nodelist*, 0 otherwise
- 2) 1 if node (MSB, LSB) in *minlist*, 0 otherwise
- 3) level of node (MSB, LSB) in *nodelist*, 0 otherwise
- 4) fanout of node (MSB, LSB) in *nodelist*, 0 otherwise

where the fanout of a node refers to the number of children it has and the level of a node refers to its topological depth from input nodes in the prefix graph. Features are normalized to $[0, 1]$.

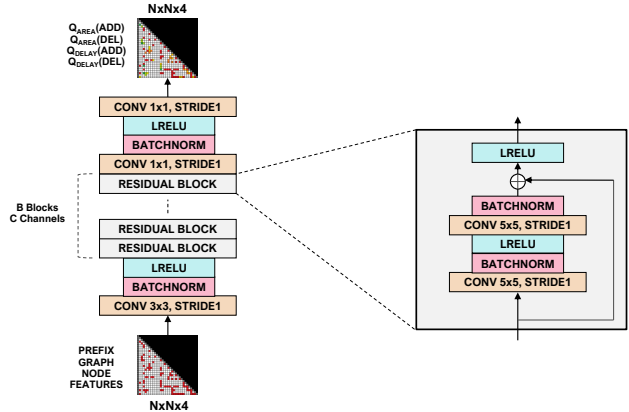


Fig. 2: Q-network architecture for N-input prefix graphs. For both 32b and 64b, our architecture uses $B = 32$ and $C = 256$.

The PrefixRL agent uses a convolutional architecture [25] in a residual network [26] configuration similar to [12]. The neural network consists of a “body” followed by a Q-value “head”, with the details of the layer types and connections given in Fig. 2. The output is a $N \times N \times 4$ matrix where the 4 channels represent the Q_{area} and Q_{delay} values for adding and for deleting the node (MSB, LSB). We use *nodelist* and *minlist* to set the Q values of illegal actions to $-\infty$ so that they are never chosen. In our experiments, we use a discount factor of $\gamma = 0.75$, an experience buffer with up to 4×10^5 elements, synchronize our target Q-network every 60 gradient steps and train using the Adam optimizer with learning rate 4×10^{-5} .

D. Training System Description

Building the deep learning training system for PrefixRL is a nontrivial task due to the time and resource-intensiveness of the reward calculation $\mathcal{R}(s, a)$. For each prefix graph state we perform circuit synthesis by generating a gate level netlist (Section V-A) and applying timing-driven synthesis optimizations at 4 delay targets. We use the OpenPhySyn physical synthesis tool [16] for optimizations such as gate sizing, gate cloning, buffer insertion and pin swapping to prefix circuits. Section V-D highlights the importance of these optimizations. After synthesis optimization, we interpolate an area-delay tradeoff curve using PChip Interpolation (Fig. 3). Processing a single state with OpenPhySyn takes up to 36 seconds on average for the 64b case (see Table I) and RL training runs usually span on the order of at least 10^5 environment steps so it is necessary to parallelize synthesis as much as possible.

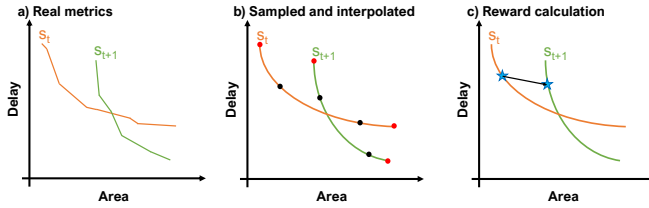


Fig. 3: Process for calculating reward in PrefixRL. (a) Each prefix state corresponds to a curve of synthesized circuit metrics based on physical synthesis timing constraints. (b) We only sample 4 such points and interpolate a curve for the metrics. (c) The reward is then calculated as the vector difference between the w -optimal points on the curves from s_t and s_{t+1} .

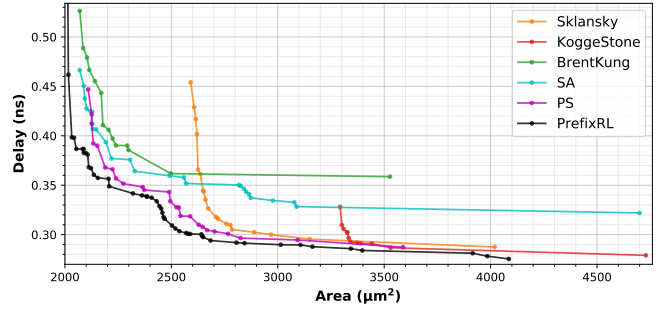
We implement an asynchronous distributed system to efficiently perform synthesis in the training loop. Our main training process launches additional processes on separate nodes with large CPU resources to run the OpenPhySyn workload and returns the final area and delay information. This asynchronous workflow is made efficient through pipeline parallelism which hides the latency of individual calls and amortizes the overall delay to be asymptotically minimal. Furthermore, we cache synthesized state designs to reduce redundant calculations and find that as the exploration parameter ϵ diminishes, the cache hit percentage becomes 50% in the 32b case and 10% in the 64b case.

Aside from parallelizing synthesis in the loop, we also built a distributed framework for efficient RL. The key observation is that DQN is an *off-policy* RL algorithm, meaning that training can be done on experience gathered by any policy, not simply the current version. This effectively decouples and allows for parallelizing the process of generating new experience and training on prior experience. For this to be viable, we must ensure an appropriate amount of new experience is generated with each new version of the parameters. We found that having 192 synthesis worker processes generating experience was empirically sufficient to satisfy this condition.

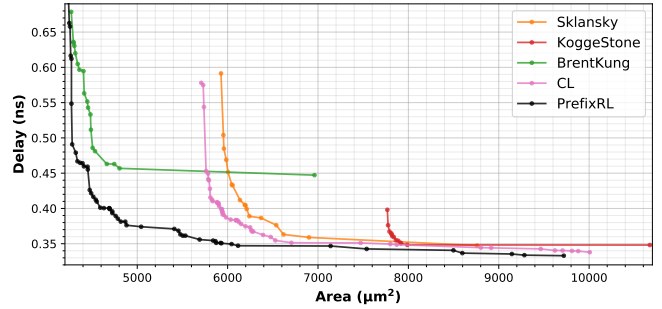
V. RESULTS AND DISCUSSION

A. Prefix Adder Optimization with OpenPhySyn

With the objective of minimizing the synthesized area and delay of prefix adder circuits, multiple PrefixRL agents were trained with 15 area-delay scalarization weights w in the range $[0.10, 0.99]$. The prefix adders were generated from prefix graphs using alternating NAND/NOR, OAI/AOI, XNOR, NOR and INV gates in the Nangate45 cell library [17] based on the implementation described in [27]. The training approach described in Section IV-D was used with the timing-driven synthesis optimizations in the OpenPhySyn [16]



(a) 32b Adder Synthesis. OpenPhySyn, Nangate45.



(b) 64b Adder Synthesis. OpenPhySyn, Nangate45.

Fig. 4: Area-delay Pareto curves for 32b and 64b adders synthesized with OpenPhySyn. PrefixRL adds Pareto-dominant all prior work. (Sklansky [3], KoggeStone [4], BrentKung [5], SA [14], PS [15], CL [10])

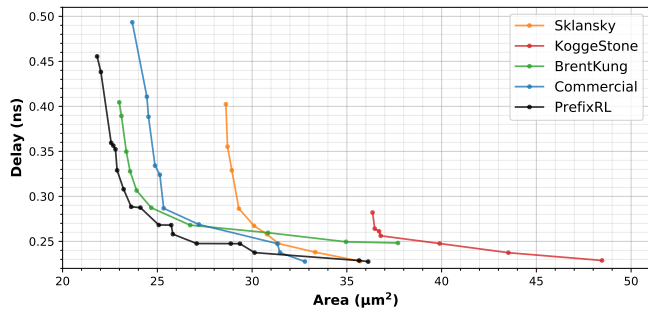
synthesis tool. Since OpenPhySyn and Nangate45 are in the open-source domain, our training procedure and results are reproducible without requiring commercial tools.

We note that circuit power is an important metric that should ideally be jointly optimized with area and delay. However, due to the computational requirements of power simulation, we did not integrate this as a third objective. We leave the integration of a power objective to the optimization as future work. Furthermore, PrefixRL agents learn policies to design prefix circuits under uniform timing constraints since the circuit synthesis environment specify uniform arrival and departure times for inputs and outputs. An interesting future direction would be to train generalizable agents that adapt to various nonuniform timing constraints.

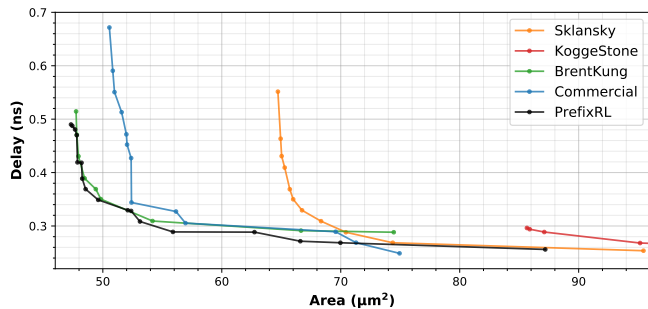
After training, the various PrefixRL agents learn to design adders specializing at various area-delay tradeoff points after considering synthesis optimizations. We synthesize the various adders generated by PrefixRL and baseline approaches at 40 delay targets with OpenPhySyn and Nangate45. Since each delay target potentially generates a different circuit for the same design, we bin all adder circuits for an approach and present the area-delay Pareto front.

The baselines for 32b prefix adders are regular Sklansky [3], Kogge-Stone [4], Brent-Kung [5] adders, and adders from simulated annealing (SA) [14] and pruned search (PS) [15] approaches. (Fig. 4a) shows that Prefix-RL agents learn to design 32b adders that Pareto dominate all these approaches. Throughout much of the delay targets (≥ 30 ns), the percent improvement in area is consistent but only 2 to 8 percentage points. However, the gains become much more significant at lower targets, reaching a maximum area saving of **16.0%** at delay target 0.293 ns.

The baselines for 64b prefix adders are the regular adders, and 1100 adders from the machine learning driven cross layer optimization approach (CL) [10]. (Fig. 4b) shows that Prefix-RL agents produce 64b adders that Pareto dominate this work as well. Samples of the learnt adders are visible in Fig. 7. Compared to the 32b setting, we



(a) 32b Adder Synthesis. Commercial tool, 8nm.



(b) 64b Adder Synthesis. Commercial tool, 8nm.

Fig. 5: Area-delay Pareto curves for 32b and 64b adders synthesized using a Commercial Synthesis Tool. Despite having been trained only with OpenPhySyn, PrefixRL-generated adders still Pareto-dominate existing baselines at all but the lowest delay target. (Sklansky [3], KoggeStone [4], BrentKung [5])

observe large gains over the baselines in the knee of the curve. This illustrates the power of an RL approach as it can successfully scale to larger problem sizes in ways that other unrestricted search space algorithms like SA cannot. In this region, PrefixRL consistently yields area savings of 12 to 20 percentage points. At lower targets, PrefixRL does even better, achieving a maximum improvement of **30.2%** at delay target 0.347 ns.

B. Generalization to Industrial Settings

To study the generalization of prefix circuit optimizations across cell libraries, we picked 7 Pareto-optimal PrefixRL adders and synthesized them at an 8nm industrial cell library with a commercial physical synthesis tool from a leading EDA vendor. In this setting, the adder is instantiated with inputs arriving from, and outputs feeding to flip-flops to ensure input uniform arrival and output departure times. We also allow the INV, XOR, XNOR cells in the adder netlist to undergo logic synthesis and technology mapping. We limit our comparison to regular adders and the library of adders instantiated by the tool (Commercial) and synthesize all the adders at the same 12 delay targets. We measure the cell area of just the adder in our results. (Fig. 5a) shows that Prefix-RL adders circuits Pareto-dominate the Kogge-Stone and Brent-Kung adders, while achieving lower area and delay than the Commercial and Sklansky adders in all the instances except the lowest delay target. This indicates that prefix circuit optimizations can generalize across cell libraries to a certain extent. We note that while training PrefixRL directly with the industrial cell library and synthesis tool may improve results further, there are possible design scenarios where PrefixRL adders may already yield better results than Commercial adders.

C. Scaling

As discussed in Section IV-D, PrefixRL requires considerable engineering optimization in order to be a tractable solution. We can

quantify this improvement by considering that 64b results were obtained after 5.0×10^5 environment steps which took approximately 5 days of training on our parallelized infrastructure. To obtain this with a single-threaded version of PrefixRL would take over 8 times longer or about 44 days of training.

TABLE I: Comparison of 16b, 32b and 64b PrefixRL adder design

Statistic	16b	32b	64b
$ \mathcal{A} $	105	465	1953
Synthesis time	11.39s	16.85s	35.56s
Train iteration time	0.45s	1.61s	3.15s
# of residual blocks	16	32	32
per-GPU batch size	96	96	6
# of data-parallel GPUs	1	1	14

Synthesis times are for Sklansky adders evaluated at 4 timing constraints. The problem space grows quickly with the number of bits and impacts other details of training.

Even with these large efficiency gains, we still had to make concessions when scaling PrefixRL to the 64b setting. The larger state representation prevents us from further expanding our 64b model capacity, so we kept it equal to that of the 32b model, while we leverage data parallelism across multiple GPUs to fit training batches in GPU memory (Table I). Training also takes roughly twice as many environment steps as needed as our 32b models to produce the results in Fig. 4a and 5a. Other common RL workloads, however, regularly reach the 10^6 to 10^7 environment step range, and in that context our solution is relatively data efficient while producing state-of-the-art results.

D. Importance of Synthesis-In-The-Loop

While we primarily focus on the physical synthesis metrics since they correspond closely to real-world performance, adder prefix graph structures are also commonly evaluated with analytical metrics. We use the analytical evaluation to assess the advantage RL provides over existing approaches as well as to examine how well performance transfers between the analytical and physical synthesis settings.

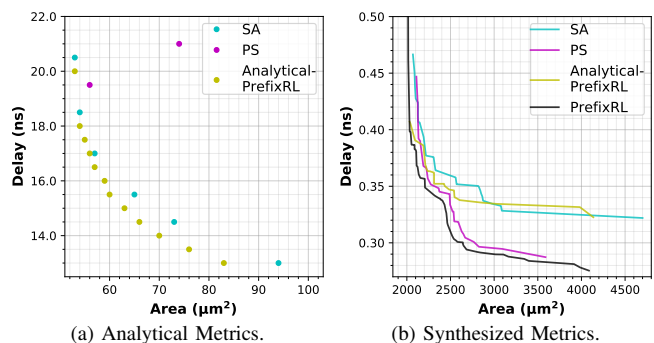


Fig. 6: 32b Analytical-PrefixRL adders outperform SA [14] and PS [15] adders when trained and compared with analytical metrics from [14]. However, timing-driven synthesis optimizations in OpenPhySyn optimize the PS adders better than the Analytical-PrefixRL and SA adders in synthesized circuits. PrefixRL adders trained with OpenPhySyn synthesis in-the-loop have the best performing synthesized circuits.

In Fig. 6a, we directly compare PrefixRL against a simulated annealing (SA) [14] approach that also operates on the unrestricted design space of all prefix graphs. We trained multiple PrefixRL agents with various area-delay scalarization weights \mathbf{w} and the same analytical models of node area (1.0) and delay ($1.0 + 0.5 \cdot \text{fanout}$) provided by [14]. The trained Analytical-PrefixRL agents learn to produce 32b prefix graphs that Pareto-dominate (Fig. 6a) all the published

solutions from [14] with **11.7%** lower area at the lowest delay point, showing that RL itself can out-compete existing approaches in a setting that does not require expensive physical synthesis feedback.

In order to compare the analytical metrics of prefix graphs against real circuit properties, we generated prefix adder circuits from the graphs and applied timing-driven synthesis optimizations at 40 target delays with OpenPhySyn and Nangate45. (Fig. 6b) shows that even though the Analytical-PrefixRL and SA designs Pareto-dominate the regular prefix graphs [3]–[5] and pruned search (PS) designs [15] at analytical metrics, they do not yield well to synthesis optimizations. The PS and Sklansky adders can achieve lower delay while maintaining lower area than the Analytical-PrefixRL and SA adders. These results highlight the importance of training directly with synthesized circuit area and delay objectives.

VI. CONCLUSION

In this paper, we have presented PrefixRL: a new deep reinforcement learning based solution for optimizing prefix circuits for synthesized area and delay. PrefixRL does not use heuristics to evaluate solutions in a pruned design space, but rather uses a deep neural network model that learns strategies purely through exploration of the unrestricted design space and feedback from synthesis tools. We apply PrefixRL to the task of designing area-delay optimized 32b and 64b prefix adders and demonstrate that it finds a frontier of designs across various area-delay trade-offs that significantly outperform previous methods. We observe that agents trained with open-source synthesis tools and cell library can design adder circuits that achieve lower area and delay than commercial tool adders in an industrial cell library. This result further suggests that using target cell libraries or commercial synthesis tools during training is a promising direction for further improvement. PrefixRL demonstrates the potential for deep reinforcement learning as an effective optimization algorithm for prefix circuits. In the future, we hope to extend the framework to other datapath circuits, consider nonuniform timing constraints and power objectives.

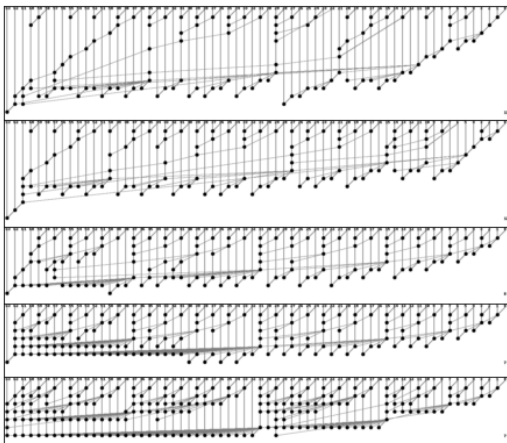


Fig. 7: 64b PrefixRL Solutions

REFERENCES

- [1] M. Lin, *Introduction to VLSI Systems: A Logic, Circuit, and System Perspective*. CRC Press, 2011.
- [2] A. K. Verma and P. Ienne, "Towards the automatic exploration of arithmetic-circuit architectures," in *Proceedings of the 43rd annual Design Automation Conference on*, 2006, pp. 445–450.
- [3] J. Sklansky, "Conditional-sum addition logic," *Ire Transactions on Electronic Computers*, vol. 9, no. 2, pp. 226–231, 1960.

- [4] P. M. Kogge and H. S. Stone, "A parallel algorithm for the efficient solution of a general class of recurrence equations," *IEEE Transactions on Computers*, vol. 22, no. 8, pp. 786–793, 1973.
- [5] Brent and Kung, "A regular layout for parallel adders," *IEEE Transactions on Computers*, vol. 31, no. 3, pp. 260–264, 1982.
- [6] T. Matsunaga and Y. Matsunaga, "Area minimization algorithm for parallel prefix adders under bitwise delay constraints," in *Proceedings of the 17th ACM Great Lakes symposium on VLSI*, 2007, pp. 435–440.
- [7] J. Liu, S. Zhou, H. Zhu, and C.-K. Cheng, "An algorithmic approach for generic parallel adders," in *ICCAD-2003. International Conference on Computer Aided Design (IEEE Cat. No.03CH37486)*, 2003, pp. 734–740.
- [8] J. P. Fishburn, "A depth-decreasing heuristic for combinational logic; or how to convert a ripple-carry adder into a carry-lookahead adder or anything in-between," in *Proceedings of the 27th ACM/IEEE Design Automation Conference on*, 1990, pp. 361–364.
- [9] R. Zimmermann, "Non-heuristic optimization and synthesis of parallel-prefix adders," in *proc. of IFIP workshop*. Citeseer, 1996.
- [10] Y. Ma, S. Roy, J. Miao, J. Chen, and B. Yu, "Cross-layer optimization for high speed adders: A pareto driven machine learning approach," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 38, no. 12, pp. 2298–2311, 2019.
- [11] R. S. Sutton and A. G. Barto, *Introduction to Reinforcement Learning*, 1st ed. Cambridge, MA, USA: MIT Press, 1998.
- [12] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. Lillicrap, K. Simonyan, and D. Hassabis, "A general reinforcement learning algorithm that masters chess, shogi, and go through self-play," *Science*, vol. 362, no. 6419, pp. 1140–1144, 2018.
- [13] A. Mirhoseini, A. Goldie, M. Yazgan, J. Jiang, E. Songhori, S. Wang, Y.-J. Lee, E. Johnson, O. Pathak, S. Bae, A. Nazi, J. Pak, A. Tong, K. Srinivasa, W. Hang, E. Tuncer, A. Babu, Q. V. Le, J. Laudon, R. Ho, R. Carpenter, and J. Dean, "Chip placement with deep reinforcement learning," 2020.
- [14] T. Moto and M. Kaneko, "Prefix sequence: Optimization of parallel prefix adders using simulated annealing," in *2018 IEEE International Symposium on Circuits and Systems (ISCAS)*, 2018, pp. 1–5.
- [15] S. Roy, M. R. Choudhury, R. Puri, and D. Z. Pan, "Towards optimal performance-area trade-off in adders by synthesis of parallel prefix structures," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 33, no. 10, pp. 1517–1530, 2014.
- [16] A. Agiza and S. Reda, "Openphysyn: An open-source physical synthesis optimization toolkit," in *2020 Workshop on Open-Source EDA Technology (WOSET)*, 2020.
- [17] Silvaco, "Nangate freepdk45 generic open cell library." [Online]. Available: <https://si2.org/open-cell-library/>
- [18] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, "Optimization by simulated annealing," *science*, vol. 220, no. 4598, pp. 671–680, 1983.
- [19] G. E. Blelloch, "Prefix sums and their applications," 1990.
- [20] R. E. Ladner and M. J. Fischer, "Parallel prefix computation," *Journal of the ACM*, vol. 27, no. 4, pp. 831–838, 1980.
- [21] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, no. 3–4, pp. 279–292, 1992.
- [22] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [23] H. van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in *AAAI'16 Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, 2016, pp. 2094–2100.
- [24] H. Mossalam, Y. M. Assael, D. M. Roijers, and S. Whiteson, "Multi-objective deep reinforcement learning," *arXiv preprint arXiv:1610.02707*, 2016.
- [25] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25*, 2012, pp. 1097–1105.
- [26] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [27] R. Zimmermann, *Binary adder architectures for cell-based VLSI and their synthesis*. Citeseer, 1997.